

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
30 March 2006 (30.03.2006)

PCT

(10) International Publication Number
WO 2006/034061 A2

(51) International Patent Classification:
C12N 15/09 (2006.01) *C12N 15/31* (2006.01)

(74) Agents: STEFFEY, Charles, E. et al.; Schwegman, Lundberg, Woessner & Kluth, P.A., P.O. Box 2938, Minneapolis, MN 55402 (US).

(21) International Application Number:
PCT/US2005/033218

(81) Designated States (*unless otherwise indicated, for every kind of national protection available*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, LY, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SM, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

(22) International Filing Date:
16 September 2005 (16.09.2005)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
10/943,508 17 September 2004 (17.09.2004) US

(84) Designated States (*unless otherwise indicated, for every kind of regional protection available*): ARIPO (BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, NL, PL, PT, RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

(71) Applicant (*for all designated States except US*):
PROMEGA CORPORATION [US/US]; 2800 Woods Hollow Road, Madison, WI 53711 (US).

(72) Inventors; and

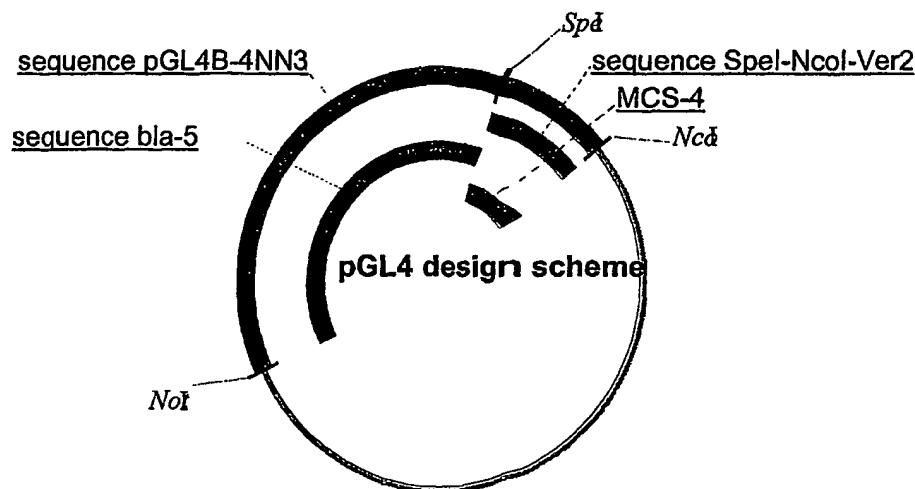
(75) Inventors/Applicants (*for US only*): WOOD, Keith, V. [US/US]; 8380 Swan Road, Mt. Horeb, WI 53572 (US). WOOD, Monika, G. [US/US]; 8380 Swan Road, Mt. Horeb, WI 53572 (US). ALMOND, Biran [US/US]; 5765 Richard Drive, Fitchburg, WI 53719 (US). PAGUIO, Aileen [US/US]; 205 Ramsey Court, Madison, WI 53704 (US). FAN, Frank [TZ/US]; 2977 Dunmore Street, Madison, WI 53711 (US).

Published:

— without international search report and to be republished upon receipt of that report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: SYNTHETIC NUCLEIC ACID MOLECULE AND METHODS OF PREPARATION



(57) Abstract: A method to prepare synthetic nucleic acid molecules having reduced inappropriate or unintended transcriptional characteristics when expressed in a particular host cell.

WO 2006/034061 A2

SYNTHETIC NUCLEIC ACID MOLECULE AND METHODS OF PREPARATION

Background

5 Transcription, the synthesis of an RNA molecule from a sequence of DNA is the first step in gene expression. Sequences which regulate DNA transcription include promoter sequences, polyadenylation signals, transcription factor binding sites and enhancer elements. A promoter is a DNA sequence capable of specific initiation of transcription and consists of three general
10 regions. The core promoter is the sequence where the RNA polymerase and its cofactors bind to the DNA. Immediately upstream of the core promoter is the proximal promoter which contains several transcription factor binding sites that are responsible for the assembly of an activation complex that in turn recruits the polymerase complex. The distal promoter, located further upstream of the
15 proximal promoter also contains transcription factor binding sites. Transcription termination and polyadenylation, like transcription initiation, are site specific and encoded by defined sequences. Enhancers are regulatory regions, containing multiple transcription factor binding sites, that can significantly increase the level of transcription from a responsive promoter regardless of the enhancer's
20 orientation and distance with respect to the promoter as long as the enhancer and promoter are located within the same DNA molecule. The amount of transcript produced from a gene may also be regulated by a post-transcriptional mechanism, the most important being RNA splicing that removes intervening sequences (introns) from a primary transcript between splice donor and splice
25 acceptor sequences.

 Natural selection is the hypothesis that genotype-environment interactions occurring at the phenotypic level lead to differential reproductive success of individuals and therefore to modification of the gene pool of a population. Some properties of nucleic acid molecules that are acted upon by
30 natural selection include codon usage frequency, RNA secondary structure, the efficiency of intron splicing, and interactions with transcription factors or other nucleic acid binding proteins. Because of the degenerate nature of the genetic

code, these properties can be optimized by natural selection without altering the corresponding amino acid sequence.

Under some conditions, it is useful to synthetically alter the natural nucleotide sequence encoding a polypeptide to better adapt the polypeptide for alternative applications. A common example is to alter the codon usage frequency of a gene when it is expressed in a foreign host cell. Although redundancy in the genetic code allows amino acids to be encoded by multiple codons, different organisms favor some codons over others. It has been found that the efficiency of protein translation in a non-native host cell can be substantially increased by adjusting the codon usage frequency but maintaining the same gene product (U.S. Patent Nos. 5,096,825, 5,670,356, and 5,874,304).

However, altering codon usage may, in turn, result in the unintentional introduction into a synthetic nucleic acid molecule of inappropriate transcription regulatory sequences. This may adversely effect transcription, resulting in anomalous expression of the synthetic DNA. Anomalous expression is defined as departure from normal or expected levels of expression. For example, transcription factor binding sites located downstream from a promoter have been demonstrated to effect promoter activity (Michael et al., 1990; Lamb et al., 1998; Johnson et al., 1998; Jones et al., 1997). Additionally, it is not uncommon for an enhancer element to exert activity and result in elevated levels of DNA transcription in the absence of a promoter sequence or for the presence of transcription regulatory sequences to increase the basal levels of gene expression in the absence of a promoter sequence.

Thus, what is needed is a method for making synthetic nucleic acid molecules with altered codon usage without also introducing inappropriate or unintended transcription regulatory sequences for expression in a particular host cell.

Summary of the Invention

The invention provides an isolated nucleic acid molecule (a polynucleotide) comprising a synthetic nucleotide sequence having reduced, for instance, 90% or less, e.g., 80%, 78%, 75%, or 70% or less, nucleic acid sequence identity relative to a parent nucleic acid sequence, e.g., a wild-type

nucleic acid sequence, and having fewer regulatory sequences such as transcription regulatory sequences. In one embodiment, the synthetic nucleotide sequence has fewer regulatory sequences than would result if the sequence differences between the synthetic nucleotide sequence and the parent nucleic acid sequence, e.g., optionally the result of differing codons, were randomly selected. In one embodiment, the synthetic nucleotide sequence encodes a polypeptide that has an amino acid sequence that is at least 85%, 90%, 95%, or 99%, or 100%, identical to the amino acid sequence of a naturally-occurring (native or wild-type) corresponding polypeptide (protein). Thus, it is recognized that some specific amino acid changes may also be desirable to alter a particular phenotypic characteristic of a polypeptide encoded by the synthetic nucleotide sequence. Preferably, the amino acid sequence identity is over at least 100 contiguous amino acid residues. In one embodiment of the invention, the codons in the synthetic nucleotide sequence that differ preferably encode the same amino acids as the corresponding codons in the parent nucleic acid sequence.

Hence, in one embodiment, the invention provides an isolated nucleic acid molecule comprising a synthetic nucleotide sequence having a coding region for a selectable or screenable polypeptide, wherein the synthetic nucleotide sequence has 90%, e.g., 80%, or less nucleic acid sequence identity to a parent nucleic acid sequence encoding a corresponding selectable or screenable polypeptide, and wherein the synthetic nucleotide sequence encodes a selectable or screenable polypeptide with at least 85% amino acid sequence identity to the corresponding selectable or screenable polypeptide encoded by the parent nucleic acid sequence. The decreased nucleotide sequence identity may be a result of different codons in the synthetic nucleotide sequence relative to the codons in the parent nucleic acid sequence. The synthetic nucleotide sequence of the invention has a reduced number of regulatory sequences relative to the parent nucleic acid sequence, for example, relative to the average number of regulatory sequences resulting from random selections of codons or nucleotides at the sequences which differ between the synthetic nucleotide sequence and the parent nucleic acid sequence. In one embodiment, a nucleic acid molecule may include a synthetic nucleotide sequence which together with other sequences encodes a selectable or screenable polypeptide. For instance, a synthetic nucleotide

sequence which forms part of an open reading frame for a selectable or screenable polypeptide may include at least 100, 150, 200, 250, 300 or more nucleotides of the open reading, which nucleotides have reduced nucleic acid sequence identity relative to corresponding sequences in a parent nucleic acid sequence. In one embodiment, the parent nucleic acid sequence is SEQ ID NO:1, SEQ ID NO:6, SEQ ID NO:15 or SEQ ID NO:41, the complement thereof, or a sequence that has 90%, 95% or 99% nucleic acid sequence identity thereto.

In one embodiment, the nucleic acid molecule of the invention comprises sequences which have been optimized for expression in mammalian cells, and more preferably, in human cells (see, e.g., WO 02/16944 which discloses methods to optimize sequences for expression in a cell of interest). For instance, nucleic acid molecules may be optimized for expression in eukaryotic cells by introducing a Kozak sequence and/or one or more introns or decreasing the number of other regulatory sequences, and/or altering codon usage to codons employed more frequently in one or more eukaryotic organisms, e.g., codons employed more frequently in an eukaryotic host cell to be transformed with the nucleic acid molecule.

In one embodiment, the synthetic nucleotide sequence is present in a vector, e.g., a plasmid, and such a vector may include other optimized sequences. In one embodiment, the synthetic nucleotide sequence encodes a polypeptide comprising a selectable polypeptide, which synthetic nucleotide sequence has at least 90% or more nucleic acid sequence identity to an open reading frame in a sequence comprising, for example, SEQ ID NO:5, SEQ ID NO:9, SEQ ID NO:10, SEQ ID NO:11, SEQ ID NO:30, SEQ ID NO:38, SEQ ID NO:39, SEQ ID NO:42, SEQ ID NO:44, SEQ ID NO:70, SEQ ID NO:71, SEQ ID NO:72, SEQ ID NO:73, SEQ ID NO:74, SEQ ID NO:80, SEQ ID NO:81, SEQ ID NO:82, SEQ ID NO:83, SEQ ID NO:84, the complement thereof, or a fragment thereof that encodes a polypeptide with substantially the same activity as the corresponding full-length and optionally wild-type (functional) polypeptide, e.g., a polypeptide encoded by SEQ ID NO:1, SEQ ID NO:6, SEQ ID NO:15 or SEQ ID NO:41, or a portion thereof which together with other parent or wild-type sequences encodes a polypeptide with substantially the same activity as the

corresponding full-length and optionally wild-type polypeptide. As used herein, “substantially the same activity” is at least about 70%, e.g., 80%, 90% or more, the activity of a corresponding full-length and optionally wild-type (functional) polypeptide. In one embodiment, an isolated nucleic acid molecule encodes a
5 fusion polypeptide comprising a selectable polypeptide.

Also provided is an isolated nucleic acid molecule comprising a synthetic nucleotide sequence having a coding region for a firefly luciferase, wherein the nucleic acid sequence identity of the synthetic nucleic acid molecule is 90% or less, e.g., 80%, 78%, 75% or less, compared to a parent nucleic acid sequence
10 encoding a firefly luciferase, e.g., a parent nucleic acid sequence having SEQ ID NO:14 or SEQ ID NO:43, which synthetic nucleotide sequence has fewer regulatory sequences including transcription regulatory sequences than would result if the sequence differences, e.g., differing codons, were randomly selected. Preferably, the synthetic nucleotide sequence encodes a polypeptide that has an
15 amino acid sequence that is at least 85%, preferably 90%, and most preferably 95% or 99% identical to the amino acid sequence of a naturally-occurring or parent polypeptide. Thus, it is recognized that some specific amino acid changes may be desirable to alter a particular phenotypic characteristic of the luciferase encoded by the synthetic nucleotide sequence. Preferably, the amino acid
20 sequence identity is over at least 100 contiguous amino acid residues. In one embodiment, the synthetic nucleotide sequence encodes a polypeptide comprising a firefly luciferase, which synthetic nucleotide sequence has at least 90% or more nucleic acid sequence identity to an open reading frame in a sequence comprising, for example, SEQ ID NO:21, SEQ ID NO:22, SEQ ID
25 NO:23, the complement thereof, or a fragment thereof that encodes a polypeptide with substantially the same activity as the corresponding full-length and optionally wild-type (functional) polypeptide, e.g., a polypeptide encoded by SEQ ID NO:14 or SEQ ID NO:43, or a portion thereof which together with other sequences encodes a firefly luciferase. For instance, a synthetic nucleotide
30 sequence which forms part of an open reading frame for a firefly luciferase may include at least 100, 150, 200, 250, 300 or more nucleotides of the open reading, which nucleotides have reduced nucleic acid sequence identity relative to corresponding sequences in a parent nucleic acid sequence.

In another embodiment, the invention provides an isolated nucleic acid molecule comprising a synthetic nucleotide sequence which does not include an open reading frame encoding a peptide or polypeptide of interest, e.g., the synthetic nucleotide sequence may have an open reading frame but it does not include sequences that encode a functional or desirable peptide or polypeptide, but may include one or more stop codons in one or more reading frames, one or more poly(A) adenylation sites, and/or a contiguous sequence for two or more restriction endonucleases (restriction enzymes), i.e., a multiple cloning region (also referred to as a multiple cloning site, "MCS"), and which is generally at least 20, e.g., at least 30, nucleotides in length and up to 1000 or more nucleotides, e.g., up to 10,000 nucleotides, which synthetic nucleotide sequence has fewer regulatory sequences such as transcription regulatory sequences relative to a corresponding parent nucleic acid sequence. In one embodiment, the synthetic nucleotide sequence which does not encode a peptide or polypeptide has 90% or less, e.g., 80%, or less nucleic acid sequence identity to a parent nucleic acid sequence, wherein the decreased sequence identity is a result of a reduced number of regulatory sequences in the synthetic nucleotide sequence relative to the parent nucleic acid sequence.

The regulatory sequences which are reduced in the synthetic nucleotide sequence include, but are not limited to, any combination of transcription factor binding sequences, intron splice sites, poly(A) adenylation sites (poly(A) sequences or poly(A) sites hereinafter), enhancer sequences, promoter modules, and/or promoter sequences, e.g., prokaryotic promoter sequences. Generally, a synthetic nucleic acid molecule lacks at least 10%, 20%, 50% or more of the regulatory sequences, for instance lacks substantially all of the regulatory sequences, e.g., 80%, 90% or more, for instance, 95% or more, of the regulatory sequences, present in a corresponding parent or wild-type nucleotide sequence. Regulatory sequences, e.g., transcription regulatory sequences, are well known in the art. The synthetic nucleotide sequence may also have a reduced number of restriction enzyme recognition sites, and may be modified to include selected sequences, e.g., sequences at or near the 5' and/or 3' ends of the synthetic nucleotide sequence such as Kozak sequences and/or desirable restriction enzyme recognition sites, for instance, restriction enzyme recognition sites useful

to introduce a synthetic nucleotide sequence to a specified location, e.g., in a multiple cloning region 5' and/or 3' to a nucleic acid sequence of interest.

In one embodiment, the synthetic nucleotide sequence of the invention has a codon composition that differs from that of the parent or wild-type nucleic acid sequence. Preferred codons for use in the invention are those which are employed more frequently than at least one other codon for the same amino acid in a particular organism and/or those that are not low-usage codons in that organism and/or those that are not low-usage codons in the organism used to clone or screen for the expression of the synthetic nucleotide sequence (for example, *E. coli*). Moreover, codons for certain amino acids (i.e., those amino acids that have three or more codons), may include two or more codons that are employed more frequently than the other (non-preferred) codon(s). The presence of codons in a synthetic nucleotide sequence that are employed more frequently in one organism than in another organism results in a synthetic nucleotide sequence which, when introduced into the cells of the organism that employs those codons more frequently, has a reduced risk of aberrant expression and/or is expressed in those cells at a level that may be greater than the expression of the wild type (unmodified) nucleic acid sequence in those cells under some conditions. For example, a synthetic nucleic acid molecule of the invention which encodes a selectable or screenable polypeptide may be expressed at a level that is greater, e.g., at least about 2, 3, 4, 5, 10-fold or more relative to that of the parent or wild-type (unmodified) nucleic acid sequence in a cell or cell extract under identical conditions (such as cell culture conditions, vector backbone, and the like). In one embodiment, the synthetic nucleotide sequence of the invention has a codon composition that differs from that of the parent or wild-type nucleic acid sequence at more than 10%, 20% or more, e.g., 30%, 35%, 40% or more than 45%, e.g., 50%, 55%, 60% or more of the codons.

In one embodiment of the invention, the codons that are different are those employed more frequently in a mammal, while in another embodiment the codons that are different are those employed more frequently in a plant. A particular type of mammal, e.g., human, may have a different set of preferred codons than another type of mammal. Likewise, a particular type of plant may have a different set of preferred codons than another type of plant. In one

embodiment of the invention, the majority of the codons which differ are ones that are preferred codons in a desired host cell and/or are not low usage codons in a particular host cell. Preferred codons for mammals (e.g., humans) and plants are known to the art (e.g., Wada et al., 1990). For example, preferred human
5 codons include, but are not limited to, CGC (Arg), CTG (Leu), AGC (Ser), ACC (Thr), CCC (Pro), GCC (Ala), GGC (Gly), GTG (Val), ACT (Ile), AAG (Lys), AAC (Asn), CAG (Gln), CAC (His), GAG (Glu), GAC (Asp), TAC (Tyr), TGC (Cys) and TTC (Phe) (Wada et al., 1990). Thus, synthetic nucleotide sequences of the invention have a codon composition which differs from a wild type
10 nucleic acid sequence by having an increased number of preferred human codons, e.g. CGC, CTG, TCT, AGC, ACC, CCC, GCC, GGC, GTG, ACT, AAG, AAC, CAG, CAC, GAG, GAC, TAC, TGC, TTC, or any combination thereof. For example, the synthetic nucleotide sequence of the invention may have an increased number of AGC serine-encoding codons, CCC proline-
15 encoding codons, and/or ACC threonine-encoding codons, or any combination thereof, relative to the parent or wild-type nucleic acid sequence. Similarly, synthetic nucleotide sequences having an increased number of codons that are employed more frequently in plants, have a codon composition which differs from a wild-type nucleic acid sequence by having an increased number of the
20 plant codons including, but not limited to, CGC (Arg), CTT (Leu), TCT (Ser), TCC (Ser), ACC (Thr), CCA (Pro), CCT (Pro), GCT (Ser), GGA (Gly), GTG (Val), ATC (Ile), ATT (Ile), AAG (Lys), AAC (Asn), CAA (Gln), CAC (His), GAG (Glu), GAC (Asp), TAC (Tyr), TGC (Cys), TTC (Phe), or any combination thereof (Murray et al., 1989). Preferred codons may differ for different types of
25 plants (Wada et al., 1990).

The nucleotide substitutions in the synthetic nucleic acid sequence may be influenced by many factors such as, for example, the desire to have an increased number of nucleotide substitutions such as those resulting in a silent nucleotide substitution (encodes the same amino acid) and/or decreased number
30 of regulatory sequences. Under some circumstances (e.g., to permit removal of a transcription factor binding site) it may be desirable to replace a non-preferred codon with a codon other than a preferred codon or a codon other than the preferred codon in order to decrease the number of regulatory sequences.

The invention also provides an expression cassette or vector. The expression cassette or vector of the invention comprises a synthetic nucleotide sequence of the invention operatively linked to a promoter that is functional in a cell or comprises a synthetic nucleotide sequence, respectively. Preferred
5 promoters are those functional in mammalian cells and those functional in plant cells. Optionally, the expression cassette may include other sequences, e.g., one or more restriction enzyme recognition sequences 5' and/or 3' to an open reading frame for a selectable polypeptide or luciferase and/or a Kozak sequence, and be a part of a larger polynucleotide molecule such as a plasmid, cosmid, artificial
10 chromosome or vector, e.g., a viral vector, which may include a multiple cloning region for other sequences, e.g., promoters, enhancers, other open reading frames and/or poly(A) sites. In one embodiment, a vector of the invention includes SEQ ID NO:88, SEQ ID NO:89, SEQ ID NO:90, the complement thereof, or a sequence which has at least 80% nucleic acid sequence identity thereto and
15 encodes a selectable and/or screenable polypeptide.

In one embodiment, the synthetic nucleotide sequence encoding a selectable or screenable polypeptide is introduced into a vector backbone, e.g., one which optionally has a poly(A) site 3' to the synthetic nucleotide sequence, a gene useful for selecting transformed prokaryotic cells which optionally is a
20 synthetic sequence, a gene useful for selecting transformed eukaryotic cells which optionally is a synthetic sequence, a noncoding region for decreasing transcription and/or translation into adjacent linked desirable open reading frames, and/or a multiple cloning region 5' and/or 3' to the synthetic nucleotide sequence encoding a selectable or screenable polypeptide which optionally
25 includes one or more protein destabilization sequences (see U.S. application Serial No. 10/664,341, filed September 16, 2003, the disclosure of which is incorporated by reference herein). In one embodiment, the vector having a synthetic nucleotide sequence encoding a selectable or screenable polypeptide may lack a promoter and/or enhancer which is operably linked to that synthetic
30 sequence. In another embodiment, the invention provides a vector comprising a promoter, e.g., a prokaryotic or eukaryotic promoter, operably linked to a synthetic nucleotide sequence encoding a selectable or screenable polypeptide. Such vectors optionally include one or more multiple cloning regions, such as

ones that are useful to introduce an additional open reading frame and/or a promoter for expression of the open reading frame which promoter optionally is different than the promoter for the selectable or screenable polypeptide, and/or a prokaryotic origin of replication. A "vector backbone" as used herein may
5 include sequences (open reading frames) useful to identify cells with those sequences, e.g., in prokaryotic cells, their promoters, an origin of replication for vector maintenance, e.g., in prokaryotic cells, and optionally one or more other sequences including multiple cloning regions e.g., for insertion of a promoter and/or open reading frame of interest, and sequences which inhibit transcription
10 and/or translation.

Also provided is a host cell comprising the synthetic nucleotide sequence of the invention, an isolated polypeptide (e.g., a fusion polypeptide encoded by the synthetic nucleotide sequence of the invention), and compositions and kits comprising the synthetic nucleotide sequence of the invention, a polypeptide
15 encoded thereby, or an expression cassette or vector comprising the synthetic nucleotide sequence in suitable container means and, optionally, instruction means. The host cell may be an eukaryotic cell such as a plant or vertebrate cell, e.g., a mammalian cell, including but not limited to a human, non-human primate, canine, feline, bovine, equine, ovine or rodent (e.g., rabbit, rat, ferret,
20 hamster, or mouse) cell or a prokaryotic cell.

The invention also provides a method to prepare a synthetic nucleotide sequence of the invention by genetically altering a parent, e.g., a wild-type or synthetic, nucleic acid sequence. The method comprises altering (e.g., decreasing or eliminating) a plurality of regulatory sequences in a parent nucleic
25 acid sequence, e.g., one which encodes a selectable or screenable polypeptide or one which does not encode a peptide or polypeptide, to yield a synthetic nucleotide sequence which has a decreased number of regulatory sequences and, if the synthetic nucleotide sequence encodes a polypeptide, it preferably encodes the same amino acids as the parent nucleic acid molecule. The transcription
30 regulatory sequences which are reduced include but are not limited to any of transcription factor binding sequences, intron splice sites, poly(A) sites, enhancer sequences, promoter modules, and/or promoter sequences. Preferably, the alteration of sequences in the synthetic nucleotide sequence does not result in an

increase in regulatory sequences. In one embodiment, the synthetic nucleotide sequence encodes a polypeptide that has at least 85%, 90%, 95% or 99%, or 100%, contiguous amino acid sequence identity to the amino acid sequence of the polypeptide encoded by the parent nucleic acid sequence.

5 Thus, in one embodiment, a method to prepare a synthetic nucleic acid molecule comprising an open reading frame is provided. The method includes altering the codons and/or regulatory sequences in a parent nucleic acid sequence which encodes a reporter protein such, as a firefly luciferase or a selectable polypeptide such as one encoding resistance to ampicillin, puromycin,
10 hygromycin or neomycin, to yield a synthetic nucleotide sequence which encodes a corresponding reporter polypeptide and which has for instance at least 10% or more, e.g., 20%, 30%, 40%, 50% or more, fewer regulatory sequences relative to the parent nucleic acid sequence. The synthetic nucleotide sequence has 90%, e.g., 85%, 80%, or 78%, or less nucleic acid sequence identity to the parent
15 nucleic acid sequence and encodes a polypeptide with at least 85% amino acid sequence identity to the polypeptide encoded by the parent nucleic acid sequence. The regulatory sequences which are altered include transcription factor binding sequences, intron splice sites, poly(A) sites, promoter modules, and/or promoter sequences. In one embodiment, the synthetic nucleic acid
20 sequence hybridizes under medium stringency hybridization but not stringent conditions to the parent nucleic acid sequence or the complement thereof. In one embodiment, the codons which differ encode the same amino acids as the corresponding codons in the parent nucleic acid sequence.

 Also provided is a synthetic (including a further synthetic) nucleotide
25 sequence prepared by the methods of the invention, e.g., a further synthetic nucleotide sequence in which introduced regulatory sequences or restriction endonuclease recognition sequences are optionally removed. Thus, the method of the invention may be employed to alter the codon usage frequency and/or decrease the number of regulatory sequences in any open reading frame or to
30 decrease the number of regulatory sequences in any nucleic acid sequence, e.g., a noncoding sequence. Preferably, the codon usage frequency in a synthetic nucleotide sequence which encodes a selectable or screenable polypeptide is altered to reflect that of the host organism desired for expression of that

nucleotide sequence while also decreasing the number of potential regulatory sequences relative to the parent nucleic acid molecule.

Also provided is a method to prepare a synthetic nucleic acid molecule which does not code for a peptide or polypeptide. The method includes

5 altering the nucleotides in a parent nucleic acid sequence having at least 20 nucleotides which optionally does not code for a functional or desirable peptide or polypeptide and which optionally may include sequences which inhibit transcription and/or translation, to yield a synthetic nucleotide sequence which does not include an open reading frame encoding a peptide or polypeptide of

10 interest, e.g., the synthetic nucleotide sequence may have an open reading frame but it does not include sequences that encode a functional or desirable peptide or polypeptide, but may include one or more stop codons in one or more reading frames, one or more poly(A) adenylation sites, and/or a contiguous sequence for two or more restriction endonucleases, i.e., a multiple cloning region. The

15 synthetic nucleotide sequence is generally at least 20, e.g., at least 30, nucleotides in length and up to 1000 or more nucleotides, e.g., up to 10,000 nucleotides, and has fewer regulatory sequences such as transcription regulatory sequences relative to a corresponding parent nucleic acid sequence which does not code for a peptide or polypeptide, e.g., a parent nucleic acid sequence which

20 optionally includes sequences which inhibit transcription and/or translation. The nucleotides are altered to reduce one or more regulatory sequences, e.g., transcription factor binding sequences, intron splice sites, poly(A) sites, enhancer sequences, promoter modules, and/or promoter sequences, in the parent nucleic acid sequence.

25 The invention also provides a method to prepare an expression vector. The method includes providing a linearized plasmid having a nucleic molecule including a synthetic nucleotide sequence of the invention which encodes a selectable or screenable polypeptide which is flanked at the 5' and/or 3' end by a multiple cloning region. The plasmid is linearized by contacting the plasmid

30 with at least one restriction endonuclease which cleaves in the multiple cloning region. The linearized plasmid and an expression cassette having ends compatible with the ends in the linearized plasmid are annealed, yielding an expression vector. In one embodiment, the plasmid is linearized by cleavage by

at least two restriction endonucleases, only one of which cleaves in the multiple cloning region.

Also provided is a method to clone a promoter or open reading frame. The method includes comprising providing a linearized plasmid having a multiple cloning region and a synthetic sequence of the invention which encodes a selectable or screenable polypeptide and/or a synthetic sequence of the invention which does not encode a peptide or polypeptide, which is plasmid is linearized by contacting the plasmid with at least two restriction endonucleases at least one of which cleaves in the multiple cloning region; and annealing the linearized plasmid with DNA having a promoter or an open reading frame with ends compatible with the ends of the linearized plasmid.

Exemplary methods to prepare synthetic sequences for firefly luciferase and a number of selectable polypeptide nucleic acid sequences, as well as non-coding regions present in a vector backbone, are described hereinbelow. For instance, the methods may produce synthetic selectable polypeptide nucleic acid molecules which exhibit similar or significantly enhanced levels of mammalian expression without negatively effecting other desirable physical or biochemical properties and which were also largely devoid of regulatory elements.

Clearly, the present invention has applications with many genes and across many fields of science including, but not limited to, life science research, agrigenetics, genetic therapy, developmental science and pharmaceutical development.

Brief Description of the Figures

- Figure 1. Codons and their corresponding amino acids.
Figure 2. Design scheme for the pGL4 vector.

Detailed Description of the Invention

Definitions

- The term "nucleic acid molecule" or "nucleic acid sequence" as used herein, refers to nucleic acid, DNA or RNA, that comprises noncoding or coding sequences. Coding sequences are necessary for the production of a polypeptide or protein precursor. The polypeptide can be encoded by a full-length coding

sequence or by any portion of the coding sequence, as long as the desired protein activity is retained. Noncoding sequences refer to nucleic acids which do not code for a polypeptide or protein precursor, and may include regulatory elements such as transcription factor binding sites, poly(A) sites, restriction endonuclease
5 sites, stop codons and/or promoter sequences.

A "synthetic" nucleic acid sequence is one which is not found in nature, i.e., it has been derived using molecular biological, chemical and/or informatic techniques.

A "nucleic acid", as used herein, is a covalently linked sequence of
10 nucleotides in which the 3' position of the pentose of one nucleotide is joined by a phosphodiester group to the 5' position of the pentose of the next, and in which the nucleotide residues (bases) are linked in specific sequence, i.e., a linear order of nucleotides. A "polynucleotide", as used herein, is a nucleic acid containing a sequence that is greater than about 100 nucleotides in length. An
15 "oligonucleotide" or "primer", as used herein, is a short polynucleotide or a portion of a polynucleotide. An oligonucleotide typically contains a sequence of about two to about one hundred bases. The word "oligo" is sometimes used in place of the word "oligonucleotide".

Nucleic acid molecules are said to have a "5'-terminus" (5' end) and a
20 "3'-terminus" (3' end) because nucleic acid phosphodiester linkages occur to the 5' carbon and 3' carbon of the pentose ring of the substituent mononucleotides. The end of a polynucleotide at which a new linkage would be to a 5' carbon is its 5' terminal nucleotide. The end of a polynucleotide at which a new linkage would be to a 3' carbon is its 3' terminal nucleotide. A terminal nucleotide, as
25 used herein, is the nucleotide at the end position of the 3'- or 5'-terminus.

DNA molecules are said to have "5' ends" and "3' ends" because mononucleotides are reacted to make oligonucleotides in a manner such that the 5' phosphate of one mononucleotide pentose ring is attached to the 3' oxygen of its neighbor in one direction via a phosphodiester linkage. Therefore, an end of
30 an oligonucleotides referred to as the "5' end" if its 5' phosphate is not linked to the 3' oxygen of a mononucleotide pentose ring and as the "3' end" if its 3' oxygen is not linked to a 5' phosphate of a subsequent mononucleotide pentose

ring.

As used herein, a nucleic acid sequence, even if internal to a larger oligonucleotide or polynucleotide, also may be said to have 5' and 3' ends. In either a linear or circular DNA molecule, discrete elements are referred to as being "upstream" or 5' of the "downstream" or 3' elements. This terminology reflects the fact that transcription proceeds in a 5' to 3' fashion along the DNA strand. Typically, promoter and enhancer elements that direct transcription of a linked gene (e.g., open reading frame or coding region) are generally located 5' or upstream of the coding region. However, enhancer elements can exert their effect even when located 3' of the promoter element and the coding region. Transcription termination and polyadenylation signals are located 3' or downstream of the coding region.

The term "codon" as used herein, is a basic genetic coding unit, consisting of a sequence of three nucleotides that specify a particular amino acid to be incorporation into a polypeptide chain, or a start or stop signal. The term "coding region" when used in reference to structural genes refers to the nucleotide sequences that encode the amino acids found in the nascent polypeptide as a result of translation of a mRNA molecule. Typically, the coding region is bounded on the 5' side by the nucleotide triplet "ATG" which encodes the initiator methionine and on the 3' side by a stop codon (e.g., TAA, TAG, TGA). In some cases the coding region is also known to initiate by a nucleotide triplet "TTG".

By "protein", "polypeptide" or "peptide" is meant any chain of amino acids, regardless of length or post-translational modification (e.g., glycosylation or phosphorylation). The nucleic acid molecules of the invention may also encode a variant of a naturally-occurring protein or a fragment thereof. Preferably, such a variant protein has an amino acid sequence that is at least 85%, preferably 90%, and most preferably 95% or 99% identical to the amino acid sequence of the naturally-occurring (native or wild-type) protein from which it is derived.

Polypeptide molecules are said to have an "amino terminus" (N-terminus) and a "carboxy terminus" (C-terminus) because peptide linkages occur between the backbone amino group of a first amino acid residue and the backbone

carboxyl group of a second amino acid residue. The terms "N-terminal" and "C-terminal" in reference to polypeptide sequences refer to regions of polypeptides including portions of the N-terminal and C-terminal regions of the polypeptide, respectively. A sequence that includes a portion of the N-terminal region of a polypeptide includes amino acids predominantly from the N-terminal half of the polypeptide chain, but is not limited to such sequences. For example, an N-terminal sequence may include an interior portion of the polypeptide sequence including bases from both the N-terminal and C-terminal halves of the polypeptide. The same applies to C-terminal regions. N-terminal and C-terminal regions may, but need not, include the amino acid defining the ultimate N-terminus and C-terminus of the polypeptide, respectively.

The term "wild-type" as used herein, refers to a gene or gene product that has the characteristics of that gene or gene product isolated from a naturally occurring source. A wild-type gene is that which is most frequently observed in a population and is thus arbitrarily designated the "wild-type" form of the gene. In contrast, the term "mutant" refers to a gene or gene product that displays modifications in sequence and/or functional properties (i.e., altered characteristics) when compared to the wild-type gene or gene product. It is noted that naturally-occurring mutants can be isolated; these are identified by the fact that they have altered characteristics when compared to the wild-type gene or gene product.

The term "recombinant protein" or "recombinant polypeptide" as used herein refers to a protein molecule expressed from a recombinant DNA molecule. In contrast, the term "native protein" is used herein to indicate a protein isolated from a naturally occurring (i.e., a nonrecombinant) source. Molecular biological techniques may be used to produce a recombinant form of a protein with identical properties as compared to the native form of the protein.

The term "fusion polypeptide" refers to a chimeric protein containing a protein of interest (e.g., luciferase) joined to a heterologous sequence (e.g., a non-luciferase amino acid or protein).

The terms "cell," "cell line," "host cell," as used herein, are used interchangeably, and all such designations include progeny or potential progeny of these designations. By "transformed cell" is meant a cell into which (or into

an ancestor of which) has been introduced a nucleic acid molecule of the invention, e.g., via transient transfection. Optionally, a nucleic acid molecule synthetic gene of the invention may be introduced into a suitable cell line so as to create a stably-transfected cell line capable of producing the protein or
5 polypeptide encoded by the synthetic gene. Vectors, cells, and methods for constructing such cell lines are well known in the art. The words "transformants" or "transformed cells" include the primary transformed cells derived from the originally transformed cell without regard to the number of transfers. All progeny may not be precisely identical in DNA content, due to
10 deliberate or inadvertent mutations. Nonetheless, mutant progeny that have the same functionality as screened for in the originally transformed cell are included in the definition of transformants.

Nucleic acids are known to contain different types of mutations. A "point" mutation refers to an alteration in the sequence of a nucleotide at a single
15 base position from the wild type sequence. Mutations may also refer to insertion or deletion of one or more bases, so that the nucleic acid sequence differs from the wild-type sequence.

The term "homology" refers to a degree of complementarity between two or more sequences. There may be partial homology or complete homology (i.e.,
20 identity). Homology is often measured using sequence analysis software (e.g., EMBOSS, the European Molecular Biology Open Software Suite available at <http://www.hgmp.mrc.ac.uk/Software/EMBOSS/overview/html>). Such software matches similar sequences by assigning degrees of homology to various substitutions, deletions, insertions, and other modifications. Conservative
25 substitutions typically include substitutions within the following groups: glycine, alanine; valine, isoleucine, leucine; aspartic acid, glutamic acid, asparagine, glutamine; serine, threonine; lysine, arginine; and phenylalanine, tyrosine.

The term "isolated" when used in relation to a nucleic acid, as in "isolated
30 oligonucleotide" or "isolated polynucleotide" refers to a nucleic acid sequence that is identified and separated from at least one contaminant with which it is ordinarily associated in its source. Thus, an isolated nucleic acid is present in a form or setting that is different from that in which it is found in nature. In

contrast, non-isolated nucleic acids (e.g., DNA and RNA) are found in the state they exist in nature. For example, a given DNA sequence (e.g., a gene) is found on the host cell chromosome in proximity to neighboring genes; RNA sequences (e.g., a specific mRNA sequence encoding a specific protein), are found in the cell as a mixture with numerous other mRNAs that encode a multitude of proteins. However, isolated nucleic acid includes, by way of example, such nucleic acid in cells ordinarily expressing that nucleic acid where the nucleic acid is in a chromosomal location different from that of natural cells, or is otherwise flanked by a different nucleic acid sequence than that found in nature.

The isolated nucleic acid or oligonucleotide may be present in single-stranded or double-stranded form. When an isolated nucleic acid or oligonucleotide is to be utilized to express a protein, the oligonucleotide contains at a minimum, the sense or coding strand (i.e., the oligonucleotide may be single-stranded), but may contain both the sense and anti-sense strands (i.e., the oligonucleotide may be double-stranded).

The term "isolated" when used in relation to a polypeptide, as in "isolated protein" or "isolated polypeptide" refers to a polypeptide that is identified and separated from at least one contaminant with which it is ordinarily associated in its source. Thus, an isolated polypeptide is present in a form or setting that is different from that in which it is found in nature. In contrast, non-isolated polypeptides (e.g., proteins and enzymes) are found in the state they exist in nature.

The term "purified" or "to purify" means the result of any process that removes some of a contaminant from the component of interest, such as a protein or nucleic acid. The percent of a purified component is thereby increased in the sample.

The term "operably linked" as used herein refer to the linkage of nucleic acid sequences in such a manner that a nucleic acid molecule capable of directing the transcription of a given gene and/or the synthesis of a desired protein molecule is produced. The term also refers to the linkage of sequences encoding amino acids in such a manner that a functional (e.g., enzymatically active, capable of binding to a binding partner, capable of inhibiting, etc.) protein or polypeptide is produced.

The term "recombinant DNA molecule" means a hybrid DNA sequence comprising at least two nucleotide sequences not normally found together in nature.

5 The term "vector" is used in reference to nucleic acid molecules into which fragments of DNA may be inserted or cloned and can be used to transfer DNA segment(s) into a cell and capable of replication in a cell. Vectors may be derived from plasmids, bacteriophages, viruses, cosmids, and the like.

The terms "recombinant vector" and "expression vector" as used herein refer to DNA or RNA sequences containing a desired coding sequence and
10 appropriate DNA or RNA sequences necessary for the expression of the operably linked coding sequence in a particular host organism. Prokaryotic expression vectors include a promoter, a ribosome binding site, an origin of replication for autonomous replication in a host cell and possibly other sequences, e.g. an optional operator sequence, optional restriction enzyme sites. A promoter is
15 defined as a DNA sequence that directs RNA polymerase to bind to DNA and to initiate RNA synthesis. Eukaryotic expression vectors include a promoter, optionally a polyadenylation signal and optionally an enhancer sequence.

A polynucleotide having a nucleotide sequence encoding a protein or polypeptide means a nucleic acid sequence comprising the coding region of a
20 gene, or in other words the nucleic acid sequence encodes a gene product. The coding region may be present in either a cDNA, genomic DNA or RNA form. When present in a DNA form, the oligonucleotide may be single-stranded (i.e., the sense strand) or double-stranded. Suitable control elements such as enhancers/promoters, splice junctions, polyadenylation signals, etc. may be
25 placed in close proximity to the coding region of the gene if needed to permit proper initiation of transcription and/or correct processing of the primary RNA transcript. Alternatively, the coding region utilized in the expression vectors of the present invention may contain endogenous enhancers/promoters, splice junctions, intervening sequences, polyadenylation signals, etc. In further
30 embodiments, the coding region may contain a combination of both endogenous and exogenous control elements.

The term "regulatory element" or "regulatory sequence" refers to a genetic element or sequence that controls some aspect of the expression of

nucleic acid sequence(s). For example, a promoter is a regulatory element that facilitates the initiation of transcription of an operably linked coding region. Other regulatory elements include, but are not limited to, transcription factor binding sites, splicing signals, polyadenylation signals, termination signals and enhancer elements.

Transcriptional control signals in eukaryotes comprise "promoter" and "enhancer" elements. Promoters and enhancers consist of short arrays of DNA sequences that interact specifically with cellular proteins involved in transcription. Promoter and enhancer elements have been isolated from a variety of eukaryotic sources including genes in yeast, insect and mammalian cells. Promoter and enhancer elements have also been isolated from viruses and analogous control elements, such as promoters, are also found in prokaryotes. The selection of a particular promoter and enhancer depends on the cell type used to express the protein of interest. Some eukaryotic promoters and enhancers have a broad host range while others are functional in a limited subset of cell types. For example, the SV40 early gene enhancer is very active in a wide variety of cell types from many mammalian species and has been widely used for the expression of proteins in mammalian cells. Two other examples of promoter/enhancer elements active in a broad range of mammalian cell types are those from the human elongation factor 1 gene (Uetsuki et al., 1989; Kim et al., 1990; and Mizushima and Nagata, 1990) and the long terminal repeats of the Rous sarcoma virus (Gorman et al., 1982); and the human cytomegalovirus (Boshart et al., 1985).

The term "promoter/enhancer" denotes a segment of DNA containing sequences capable of providing both promoter and enhancer functions (i.e., the functions provided by a promoter element and an enhancer element as described above). For example, the long terminal repeats of retroviruses contain both promoter and enhancer functions. The enhancer/promoter may be "endogenous" or "exogenous" or "heterologous." An "endogenous" enhancer/promoter is one that is naturally linked with a given gene in the genome. An "exogenous" or "heterologous" enhancer/promoter is one that is placed in juxtaposition to a gene by means of genetic manipulation (i.e., molecular biological techniques) such that transcription of the gene is directed by the linked enhancer/promoter.

The presence of "splicing signals" on an expression vector often results in higher levels of expression of the recombinant transcript in eukaryotic host cells.

Splicing signals mediate the removal of introns from the primary RNA transcript and consist of a splice donor and acceptor site (Sambrook et al., 1989).

5 A commonly used splice donor and acceptor site is the splice junction from the 16S RNA of SV40.

Efficient expression of recombinant DNA sequences in eukaryotic cells requires expression of signals directing the efficient termination and polyadenylation of the resulting transcript. Transcription termination signals are
10 generally found downstream of the polyadenylation signal and are a few hundred nucleotides in length. The term "poly(A) site" or "poly(A) sequence" as used herein denotes a DNA sequence which directs both the termination and polyadenylation of the nascent RNA transcript. Efficient polyadenylation of the recombinant transcript is desirable, as transcripts lacking a poly(A) tail are
15 unstable and are rapidly degraded. The poly(A) signal utilized in an expression vector may be "heterologous" or "endogenous." An endogenous poly(A) signal is one that is found naturally at the 3' end of the coding region of a given gene in the genome. A heterologous poly(A) signal is one which has been isolated from one gene and positioned 3' to another gene. A commonly used heterologous
20 poly(A) signal is the SV40 poly(A) signal. The SV40 poly(A) signal is contained on a 237 bp *Bam*H I/*Bcl* I restriction fragment and directs both termination and polyadenylation (Sambrook et al., 1989).

Eukaryotic expression vectors may also contain "viral replicons" or "viral origins of replication." Viral replicons are viral DNA sequences which allow for
25 the extrachromosomal replication of a vector in a host cell expressing the appropriate replication factors. Vectors containing either the SV40 or polyoma virus origin of replication replicate to high copy number (up to 10^4 copies/cell) in cells that express the appropriate viral T antigen. In contrast, vectors containing the replicons from bovine papillomavirus or Epstein-Barr virus replicate
30 extrachromosomally at low copy number (about 100 copies/cell).

The term "*in vitro*" refers to an artificial environment and to processes or reactions that occur within an artificial environment. *In vitro* environments include, but are not limited to, test tubes and cell lysates. The term "*in vivo*"

refers to the natural environment (e.g., an animal or a cell) and to processes or reactions that occur within a natural environment.

The term "expression system" refers to any assay or system for determining (e.g., detecting) the expression of a gene of interest. Those skilled
 5 in the field of molecular biology will understand that any of a wide variety of expression systems may be used. A wide range of suitable mammalian cells are available from a wide range of sources (e.g., the American Type Culture Collection, Rockland, MD). The method of transformation or transfection and the choice of expression vehicle will depend on the host system selected.

10 Transformation and transfection methods are described, e.g., in Ausubel et al., 1992. Expression systems include *in vitro* gene expression assays where a gene of interest (e.g., a reporter gene) is linked to a regulatory sequence and the expression of the gene is monitored following treatment with an agent that inhibits or induces expression of the gene. Detection of gene expression can be
 15 through any suitable means including, but not limited to, detection of expressed mRNA or protein (e.g., a detectable product of a reporter gene) or through a detectable change in the phenotype of a cell expressing the gene of interest. Expression systems may also comprise assays where a cleavage event or other nucleic acid or cellular change is detected.

20 All amino acid residues identified herein are in the natural L-configuration. In keeping with standard polypeptide nomenclature, abbreviations for amino acid residues are as shown in the following Table of Correspondence.

25

| TABLE OF CORRESPONDENCE | | |
|-------------------------|----------|-----------------|
| 1-Letter | 3-Letter | AMINO ACID |
| Y | Tyr | L-tyrosine |
| G | Gly | L-glycine |
| F | Phe | L-phenylalanine |
| 30 M | Met | L-methionine |
| A | Ala | L-alanine |
| S | Ser | L-serine |
| I | Ile | L-isoleucine |

| | | | |
|----|---|-----|-----------------|
| | L | Leu | L-leucine |
| | T | Thr | L-threonine |
| | V | Val | L-valine |
| | P | Pro | L-proline |
| 5 | K | Lys | L-lysine |
| | H | His | L-histidine |
| | Q | Gln | L-glutamine |
| | E | Glu | L-glutamic acid |
| | W | Trp | L-tryptophan |
| 10 | R | Arg | L-arginine |
| | D | Asp | L-aspartic acid |
| | N | Asn | L-asparagine |
| | C | Cys | L-cysteine |

The terms "complementary" or "complementarity" are used in reference
 15 to a sequence of nucleotides related by the base-pairing rules. For example, for
 the sequence 5' "A-G-T" 3', is complementary to the sequence 3' "T-C-A" 5'.
 Complementarity may be "partial," in which only some of the nucleic acids'
 bases are matched according to the base pairing rules. Or, there may be
 "complete" or "total" complementarity between the nucleic acids. The degree of
 20 complementarity between nucleic acid strands has significant effects on the
 efficiency and strength of hybridization between nucleic acid strands. This is of
 particular importance in amplification reactions, as well as detection methods
 which depend upon hybridization of nucleic acids.

When used in reference to a double-stranded nucleic acid sequence such
 25 as a cDNA or a genomic clone, the term "substantially homologous" refers to any
 probe which can hybridize to either or both strands of the double-stranded
 nucleic acid sequence under conditions of low stringency as described herein.

"Probe" refers to an oligonucleotide designed to be sufficiently
 complementary to a sequence in a denatured nucleic acid to be probed (in
 30 relation to its length) and is bound under selected stringency conditions.

"Hybridization" and "binding" in the context of probes and denatured
 nucleic acids are used interchangeably. Probes that are hybridized or bound to
 denatured nucleic acids are base paired to complementary sequences in the

polynucleotide. Whether or not a particular probe remains base paired with the polynucleotide depends on the degree of complementarity, the length of the probe, and the stringency of the binding conditions. The higher the stringency, the higher must be the degree of complementarity and/or the longer the probe.

5 The term "hybridization" is used in reference to the pairing of complementary nucleic acid strands. Hybridization and the strength of hybridization (i.e., the strength of the association between nucleic acid strands) is impacted by many factors well known in the art including the degree of complementarity between the nucleic acids, stringency of the conditions
10 involved such as the concentration of salts, the T_m (melting temperature) of the formed hybrid, the presence of other components (e.g., the presence or absence of polyethylene glycol), the molarity of the hybridizing strands and the G:C content of the nucleic acid strands.

 The term "stringency" is used in reference to the conditions of
15 temperature, ionic strength, and the presence of other compounds, under which nucleic acid hybridizations are conducted. With "high stringency" conditions, nucleic acid base pairing will occur only between nucleic acid fragments that have a high frequency of complementary base sequences. Thus, conditions of "medium" or "low" stringency are often required when it is desired that nucleic
20 acids that are not completely complementary to one another be hybridized or annealed together. The art knows well that numerous equivalent conditions can be employed to comprise medium or low stringency conditions. The choice of hybridization conditions is generally evident to one skilled in the art and is usually guided by the purpose of the hybridization, the type of hybridization
25 (DNA-DNA or DNA-RNA), and the level of desired relatedness between the sequences (e.g., Sambrook et al., 1989; Nucleic Acid Hybridization, A Practical Approach, IRL Press, Washington D.C., 1985, for a general discussion of the methods).

 The stability of nucleic acid duplexes is known to decrease with
30 increasing numbers of mismatched bases, and further to be decreased to a greater or lesser degree depending on the relative positions of mismatches in the hybrid duplexes. Thus, the stringency of hybridization can be used to maximize or minimize stability of such duplexes. Hybridization stringency can be altered by:

adjusting the temperature of hybridization; adjusting the percentage of helix destabilizing agents, such as formamide, in the hybridization mix; and adjusting the temperature and/or salt concentration of the wash solutions. For filter hybridizations, the final stringency of hybridizations often is determined by the salt concentration and/or temperature used for the post-hybridization washes.

"High stringency conditions" when used in reference to nucleic acid hybridization comprise conditions equivalent to binding or hybridization at 42°C in a solution consisting of 5X SSPE (43.8 g/l NaCl, 6.9 g/l NaH₂PO₄ H₂O and 1.85 g/l EDTA, pH adjusted to 7.4 with NaOH), 0.5% SDS, 5X Denhardt's reagent and 100 µg/ml denatured salmon sperm DNA followed by washing in a solution comprising 0.1X SSPE, 1.0% SDS at 42°C when a probe of about 500 nucleotides in length is employed.

"Medium stringency conditions" when used in reference to nucleic acid hybridization comprise conditions equivalent to binding or hybridization at 42°C in a solution consisting of 5X SSPE (43.8 g/l NaCl, 6.9 g/l NaH₂PO₄ H₂O and 1.85 g/l EDTA, pH adjusted to 7.4 with NaOH), 0.5% SDS, 5X Denhardt's reagent and 100 µg/ml denatured salmon sperm DNA followed by washing in a solution comprising 1.0X SSPE, 1.0% SDS at 42°C when a probe of about 500 nucleotides in length is employed.

"Low stringency conditions" comprise conditions equivalent to binding or hybridization at 42°C in a solution consisting of 5X SSPE (43.8 g/l NaCl, 6.9 g/l NaH₂PO₄ H₂O and 1.85 g/l EDTA, pH adjusted to 7.4 with NaOH), 0.1% SDS, 5X Denhardt's reagent [50X Denhardt's contains per 500 ml: 5 g Ficoll (Type 400, Pharmacia), 5 g BSA (Fraction V; Sigma)] and 100 g/ml denatured salmon sperm DNA followed by washing in a solution comprising 5X SSPE, 0.1% SDS at 42°C when a probe of about 500 nucleotides in length is employed.

The term "T_m" is used in reference to the "melting temperature". The melting temperature is the temperature at which 50% of a population of double-stranded nucleic acid molecules becomes dissociated into single strands. The equation for calculating the T_m of nucleic acids is well-known in the art. The T_m of a hybrid nucleic acid is often estimated using a formula adopted from hybridization assays in 1 M salt, and commonly used for calculating T_m for PCR primers: [(number of A + T) x 2°C + (number of G+C) x 4°C]. (C.R. Newton et

al., PCR, 2nd Ed., Springer-Verlag (New York, 1997), p. 24). This formula was found to be inaccurate for primers longer than 20 nucleotides. (Id.) Another simple estimate of the T_m value may be calculated by the equation: $T_m = 81.5 + 0.41(\% G + C)$, when a nucleic acid is in aqueous solution at 1 M NaCl. (e.g.,
5 Anderson and Young, Quantitative Filter Hybridization, in Nucleic Acid Hybridization, 1985). Other more sophisticated computations exist in the art which take structural as well as sequence characteristics into account for the calculation of T_m . A calculated T_m is merely an estimate; the optimum temperature is commonly determined empirically.

10 The term "promoter/enhancer" denotes a segment of DNA containing sequences capable of providing both promoter and enhancer functions (i.e., the functions provided by a promoter element and an enhancer element as described above). For example, the long terminal repeats of retroviruses contain both promoter and enhancer functions. The enhancer/promoter may be "endogenous"
15 or "exogenous" or "heterologous." An "endogenous" enhancer/promoter is one that is naturally linked with a given gene in the genome. An "exogenous" or "heterologous" enhancer/promoter is one that is placed in juxtaposition to a gene by means of genetic manipulation (i.e., molecular biological techniques) such that transcription of the gene is directed by the linked enhancer/promoter.

20 The term "sequence homology" means the proportion of base matches between two nucleic acid sequences or the proportion of amino acid matches between two amino acid sequences. When sequence homology is expressed as a percentage, e.g., 50%, the percentage denotes the proportion of matches over the length of sequence from one sequence that is compared to some other sequence.
25 Gaps (in either of the two sequences) are permitted to maximize matching; gap lengths of 15 bases or less are usually used, 6 bases or less are preferred with 2 bases or less more preferred. When using oligonucleotides as probes or treatments, the sequence homology between the target nucleic acid and the oligonucleotide sequence is generally not less than 17 target base matches out of
30 20 possible oligonucleotide base pair matches (85%); preferably not less than 9 matches out of 10 possible base pair matches (90%), and more preferably not less than 19 matches out of 20 possible base pair matches (95%).

Two amino acid sequences are homologous if there is a partial or complete identity between their sequences. For example, 85% homology means that 85% of the amino acids are identical when the two sequences are aligned for maximum matching. Gaps (in either of the two sequences being matched) are
5 allowed in maximizing matching; gap lengths of 5 or less are preferred with 2 or less being more preferred. Alternatively and preferably, two protein sequences (or polypeptide sequences derived from them of at least 100 amino acids in length) are homologous, as this term is used herein, if they have an alignment score of at more than 5 (in standard deviation units) using the program ALIGN
10 with the mutation data matrix and a gap penalty of 6 or greater. See Dayhoff, M. O., in Atlas of Protein Sequence and Structure, 1972, volume 5, National Biomedical Research Foundation, pp. 101-110, and Supplement 2 to this volume, pp. 1-10. The two sequences or parts thereof are more preferably homologous if their amino acids are greater than or equal to 85% identical when
15 optimally aligned using the ALIGN program.

The following terms are used to describe the sequence relationships between two or more polynucleotides: "reference sequence", "comparison window", "sequence identity", "percentage of sequence identity", and "substantial identity". A "reference sequence" is a defined sequence used as a
20 basis for a sequence comparison; a reference sequence may be a subset of a larger sequence, for example, as a segment of a full-length cDNA or gene sequence given in a sequence listing, or may comprise a complete cDNA or gene sequence. Generally, a reference sequence is at least 20 nucleotides in length, frequently at least 25 nucleotides in length, and often at least 50 or 100
25 nucleotides in length. Since two polynucleotides may each (1) comprise a sequence (i.e., a portion of the complete polynucleotide sequence) that is similar between the two polynucleotides, and (2) may further comprise a sequence that is divergent between the two polynucleotides, sequence comparisons between two (or more) polynucleotides are typically performed by comparing sequences of the
30 two polynucleotides over a "comparison window" to identify and compare local regions of sequence similarity.

A "comparison window", as used herein, refers to a conceptual segment of at least 20 contiguous nucleotides and wherein the portion of the

polynucleotide sequence in the comparison window may comprise additions or deletions (i.e., gaps) of 20 percent or less as compared to the reference sequence (which does not comprise additions or deletions) for optimal alignment of the two sequences.

5 Methods of alignment of sequences for comparison are well known in the art. Thus, the determination of percent identity between any two sequences can be accomplished using a mathematical algorithm. Preferred, non-limiting examples of such mathematical algorithms are the algorithm of Myers and Miller (1988); the local homology algorithm of Smith and Waterman (1981); the
10 homology alignment algorithm of Needleman and Wunsch (1970); the search-for-similarity-method of Pearson and Lipman (1988); the algorithm of Karlin and Altschul (1990), modified as in Karlin and Altschul (1993).

Computer implementations of these mathematical algorithms can be utilized for comparison of sequences to determine sequence identity. Such
15 implementations include, but are not limited to: ClustalW (available, e.g., at <http://www.ebi.ac.uk/clustalw/>); the ALIGN program (Version 2.0) and GAP, BESTFIT, BLAST, FASTA, and TFASTA in the Wisconsin Genetics Software Package, Version 8. Alignments using these programs can be performed using the default parameters. The CLUSTAL program is well described by Higgins et
20 al. (1988); Higgins et al. (1989); Corpet et al. (1988); Huang et al. (1992); and Pearson et al. (1994). The ALIGN program is based on the algorithm of Myers and Miller, *supra*. The BLAST programs of Altschul et al. (1990), are based on the algorithm of Karlin and Altschul *supra*. To obtain gapped alignments for comparison purposes, Gapped BLAST (in BLAST 2.0) can be utilized as
25 described in Altschul et al. (1997). Alternatively, PSI-BLAST (in BLAST 2.0) can be used to perform an iterated search that detects distant relationships between molecules. See Altschul et al., *supra*. When utilizing BLAST, Gapped BLAST, PSI-BLAST, the default parameters of the respective programs (e.g. BLASTN for nucleotide sequences, BLASTX for proteins) can be used. See
30 <http://www.ncbi.nlm.nih.gov>. Alignment may also be performed manually by inspection

The term "sequence identity" means that two polynucleotide sequences are identical (i.e., on a nucleotide-by-nucleotide basis) over the window of

comparison. The term "percentage of sequence identity" means that two polynucleotide sequences are identical (i.e., on a nucleotide-by-nucleotide basis) for the stated proportion of nucleotides over the window of comparison. The term "percentage of sequence identity" is calculated by comparing two optimally aligned sequences over the window of comparison, determining the number of positions at which the identical nucleic acid base (e.g., A, T, C, G, U, or I) occurs in both sequences to yield the number of matched positions, dividing the number of matched positions by the total number of positions in the window of comparison (i.e., the window size), and multiplying the result by 100 to yield the percentage of sequence identity. The terms "substantial identity" as used herein denote a characteristic of a polynucleotide sequence, wherein the polynucleotide comprises a sequence that has at least 60%, preferably at least 65%, more preferably at least 70%, up to about 85%, and even more preferably at least 90 to 95%, more usually at least 99%, sequence identity as compared to a reference sequence over a comparison window of at least 20 nucleotide positions, frequently over a window of at least 20-50 nucleotides, and preferably at least 300 nucleotides, wherein the percentage of sequence identity is calculated by comparing the reference sequence to the polynucleotide sequence which may include deletions or additions which total 20 percent or less of the reference sequence over the window of comparison. The reference sequence may be a subset of a larger sequence.

As applied to polypeptides, the term "substantial identity" means that two peptide sequences, when optimally aligned, such as by the programs GAP or BESTFIT using default gap weights, share at least about 85% sequence identity, preferably at least about 90% sequence identity, more preferably at least about 95 % sequence identity, and most preferably at least about 99 % sequence identity.

Synthetic Nucleotide Sequences and Methods of the Invention

The invention provides compositions comprising synthetic nucleotide sequences, as well as methods for preparing those sequences which yield synthetic nucleotide sequences that are efficiently expressed as a polypeptide or protein with desirable characteristics including reduced inappropriate or

unintended transcription characteristics, or do not result in inappropriate or unintended transcription characteristics, when present in a particular cell type.

Natural selection is the hypothesis that genotype-environment interactions occurring at the phenotypic level lead to differential reproductive success of individuals and hence to modification of the gene pool of a population. It is generally accepted that the amino acid sequence of a protein found in nature has undergone optimization by natural selection. However, amino acids exist within the sequence of a protein that do not contribute significantly to the activity of the protein and these amino acids can be changed to other amino acids with little or no consequence. Furthermore, a protein may be useful outside its natural environment or for purposes that differ from the conditions of its natural selection. In these circumstances, the amino acid sequence can be synthetically altered to better adapt the protein for its utility in various applications.

Likewise, the nucleic acid sequence that encodes a protein is also optimized by natural selection. The relationship between coding DNA and its transcribed RNA is such that any change to the DNA affects the resulting RNA. Thus, natural selection works on both molecules simultaneously. However, this relationship does not exist between nucleic acids and proteins. Because multiple codons encode the same amino acid, many different nucleotide sequences can encode an identical protein. A specific protein composed of 500 amino acids can theoretically be encoded by more than 10^{150} different nucleic acid sequences.

Natural selection acts on nucleic acids to achieve proper encoding of the corresponding protein. Presumably, other properties of nucleic acid molecules are also acted upon by natural selection. These properties include codon usage frequency, RNA secondary structure, the efficiency of intron splicing, and interactions with transcription factors or other nucleic acid binding proteins. These other properties may alter the efficiency of protein translation and the resulting phenotype. Because of the redundant nature of the genetic code, these other attributes can be optimized by natural selection without altering the corresponding amino acid sequence.

Under some conditions, it is useful to synthetically alter the natural nucleotide sequence encoding a protein to better adapt the protein for alternative

applications. A common example is to alter the codon usage frequency of a gene when it is expressed in a foreign host. Although redundancy in the genetic code allows amino acids to be encoded by multiple codons, different organisms favor some codons over others. The codon usage frequencies tend to differ most for
5 organisms with widely separated evolutionary histories. It has been found that when transferring genes between evolutionarily distant organisms, the efficiency of protein translation can be substantially increased by adjusting the codon usage frequency (see U.S. Patent Nos. 5,096,825, 5,670,356 and 5,874,304).

In one embodiment, the sequence of a reporter gene is modified as the
10 codon usage of reporter genes often does not correspond to the optimal codon usage of the experimental cells. In another embodiment, the sequence of a reporter gene is modified to remove regulatory sequences such as those which may alter expression of the reporter gene or a linked gene. Examples include β -galactosidase (β -gal) and chloramphenicol acetyltransferase (*cat*) reporter genes
15 that are derived from *E. coli* and are commonly used in mammalian cells; the β -glucuronidase (*gus*) reporter gene that is derived from *E. coli* and commonly used in plant cells; the firefly luciferase (*luc*) reporter gene that is derived from an insect and commonly used in plant and mammalian cells; and the *Renilla* luciferase, and green fluorescent protein (*gfp*) reporter genes which are derived
20 from coelenterates and are commonly used in plant and mammalian cells. To achieve sensitive quantitation of reporter gene expression, the activity of the gene product must not be endogenous to the experimental host cells. Thus, reporter genes are usually selected from organisms having unique and distinctive phenotypes. Consequently, these organisms often have widely separated
25 evolutionary histories from the experimental host cells.

Previously, to create genes having a more optimal codon usage frequency but still encoding the same gene product, a synthetic nucleic acid sequence was made by replacing existing codons with codons that were generally more favorable to the experimental host cell (see U.S. Patent Nos. 5,096,825,
30 5,670,356 and 5,874,304.) The result was a net improvement in codon usage frequency of the synthetic gene. However, the optimization of other attributes was not considered and so these synthetic genes likely did not reflect genes optimized by natural selection.

In particular, improvements in codon usage frequency are intended only for optimization of a RNA sequence based on its role in translation into a protein. Thus, previously described methods did not address how the sequence of a synthetic gene affects the role of DNA in transcription into RNA. Most notably, consideration had not been given as to how transcription factors may interact with the synthetic DNA and consequently modulate or otherwise influence gene transcription. For genes found in nature, the DNA would be optimally transcribed by the native host cell and would yield an RNA that encodes a properly folded gene product. In contrast, synthetic genes have previously not been optimized for transcriptional characteristics. Rather, this property has been ignored or left to chance.

This concern is important for all genes, but particularly important for reporter genes, which are most commonly used to quantitate transcriptional behavior in the experimental host cells, and vector backbone sequences for genes. Hundreds of transcription factors have been identified in different cell types under different physiological conditions, and likely more exist but have not yet been identified. All of these transcription factors can influence the transcription of an introduced gene or sequences linked thereto. A useful synthetic reporter gene or vector backbone of the invention has a minimal risk of influencing or perturbing intrinsic transcriptional characteristics of the host cell because the structure of that gene or vector backbone has been altered. A particularly useful synthetic reporter gene or vector backbone will have desirable characteristics under a new set and/or a wide variety of experimental conditions. To best achieve these characteristics, the structure of the synthetic gene or synthetic vector backbone should have minimal potential for interacting with transcription factors within a broad range of host cells and physiological conditions. Minimizing potential interactions between a reporter gene or vector backbone and a host cell's endogenous transcription factors increases the value of a reporter gene or vector backbone by reducing the risk of inappropriate transcriptional characteristics of the gene or vector backbone within a particular experiment, increasing applicability of the gene or vector backbone in various environments, and increasing the acceptance of the resulting experimental data.

In contrast, a reporter gene comprising a native nucleotide sequence, based on a genomic or cDNA clone from the original host organism, or a vector backbone comprising native sequences found in one or a variety of different organisms, may interact with transcription factors when present in an exogenous host. This risk stems from two circumstances. First, the native nucleotide sequence contains sequences that were optimized through natural selection to influence gene transcription within the native host organism. However, these sequences might also influence transcription when the sequences are present in exogenous hosts, i.e., out of context, thus interfering with its performance as a reporter gene or vector backbone. Second, the nucleotide sequence may inadvertently interact with transcription factors that were not present in the native host organism, and thus did not participate in its natural selection. The probability of such inadvertent interactions increases with greater evolutionary separation between the experimental cells and the native organism of the reporter gene or vector backbone.

These potential interactions with transcription factors would likely be disrupted when using a synthetic reporter gene having alterations in codon usage frequency. However, a synthetic reporter gene sequence, designed by choosing codons based only on codon usage frequency, or randomly replacing sequences or randomly juxtaposing sequences in a vector backbone, is likely to contain other unintended transcription factor binding sites since the resulting sequence has not been subjected to the benefit of natural selection to correct inappropriate transcriptional activities. Inadvertent interactions with transcription factors could also occur whenever an encoded amino acid sequence is artificially altered, e.g., to introduce amino acid substitutions. Similarly, these changes have not been subjected to natural selection, and thus may exhibit undesired characteristics.

Thus, the invention provides a method for preparing synthetic nucleotide sequences that reduce the risk of undesirable interactions of the nucleotide sequence with transcription factors and other trans-acting factors when expressed in a particular host cell, thereby reducing inappropriate or unintended characteristics. Preferably, the method yields synthetic genes containing improved codon usage frequencies for a particular host cell and with a reduced

occurrence of regulatory sequences such as transcription factor binding sites and/or vector backbone sequences with a reduced occurrence of regulatory sequences. The invention also provides a method of preparing synthetic genes containing improved codon usage frequencies with a reduced occurrence of
5 transcription factor binding sites and additional beneficial structural attributes. Such additional attributes include the absence of inappropriate RNA splicing junctions, poly(A) addition signals, undesirable restriction enzyme recognition sites, ribosomal binding sites, and/or secondary structural motifs such as hairpin loops.

10 In one embodiment, a parent nucleic acid sequence encoding a polypeptide is optimized for expression in a particular cell. For example, the nucleic acid sequence is optimized by replacing codons in the wild-type sequence with codons which are preferentially employed in a particular (selected) cell, which codon replacement also reduces the number of regulatory
15 sequences. Preferred codons have a relatively high codon usage frequency in a selected cell, and preferably their introduction results in the introduction of relatively few regulatory sequences such as transcription factor binding sites, and relatively few other undesirable structural attributes. Thus, the optimized nucleotide sequence may have an improved level of expression due to improved
20 codon usage frequency, and a reduced risk of inappropriate transcriptional behavior due to a reduced number of undesirable transcription regulatory sequences. In another embodiment, a parent vector backbone sequence is altered to remove regulatory sequences and optionally restriction endonuclease sites, and optionally retain or add other desirable characteristics, e.g., the presence of one
25 or more stop codons in one or more reading frames, one or more poly(A) sites, and/or restriction endonuclease sites.

The invention may be employed with any nucleic acid sequence, e.g., a native sequence such as a cDNA or one that has been manipulated *in vitro*. Exemplary genes include, but are not limited to, those encoding lactamase (β -gal), neomycin resistance (Neo), hygromycin resistance (Hyg), puromycin
30 resistance (Puro), ampicillin resistance (Amp), CAT, GUS, galactopyranoside, GFP, xylosidase, thymidine kinase, arabinosidase, luciferase and the like. As used herein, a "reporter gene" is a gene that imparts a distinct phenotype to cells

expressing the gene and thus permits cells having the gene to be distinguished from cells that do not have the gene. Such genes may encode either a selectable or screenable polypeptide, depending on whether the marker confers a trait which one can 'select' for by chemical means, i.e., through the use of a selective agent (e.g., a herbicide, antibiotic, or the like), or whether it is simply a "reporter" trait that one can identify through observation or testing, i.e., by 'screening'. Included within the terms selectable or screenable marker genes are also genes which encode a "secretable marker" whose secretion can be detected as a means of identifying or selecting for transformed cells. Examples include markers that encode a secretable antigen that can be identified by antibody interaction, or even secretable enzymes which can be detected by their catalytic activity. Secretable proteins fall into a number of classes, including small, diffusible proteins detectable, e.g., by ELISA, and proteins that are inserted or trapped in the cell membrane.

Elements of the present disclosure are exemplified in detail through the use of particular genes and vector backbone sequences. Of course, many examples of suitable genes and vector backbones are known to the art and can be employed in the practice of the invention. Therefore, it will be understood that the following discussion is exemplary rather than exhaustive. In light of the techniques disclosed herein and the general recombinant techniques that are known in the art, the present invention renders possible the alteration of any gene or vector backbone sequence.

Exemplary genes include, but are not limited to, a *neo* gene, a *puro* gene, an *amp* gene, a β -gal gene, a *gus* gene, a *cat* gene, a *gpt* gene, a *hyg* gene, a *hisD* gene, a *ble* gene, a *mppt* gene, a *bar* gene, a nitrilase gene, a mutant acetolactate synthase gene (ALS) or acetoacid synthase gene (AAS), a methotrexate-resistant *dhfr* gene, a dalapon dehalogenase gene, a mutated anthranilate synthase gene that confers resistance to 5-methyl tryptophan (WO 97/26366), an R-locus gene, a β -lactamase gene, a *xyIE* gene, an α -amylase gene, a tyrosinase gene, a luciferase (*luc*) gene (e.g., a *Renilla reniformis* luciferase gene, a firefly luciferase gene, or a click beetle luciferase (*Pyrophorus plagiophthalmus* gene), an aequorin gene, or a fluorescent protein gene.

The method of the invention can be performed by, although it is not limited to, a recursive process. The process includes assigning preferred codons to each amino acid in a target molecule, e.g., a native nucleotide sequence, based on codon usage in a particular species, identifying potential transcription regulatory sequences such as transcription factor binding sites in the nucleic acid sequence having preferred codons, e.g., using a database of such binding sites, optionally identifying other undesirable sequences, and substituting an alternative codon (i.e., encoding the same amino acid) at positions where undesirable transcription factor binding sites or other sequences occur. For codon distinct versions, alternative preferred codons are substituted in each version. If necessary, the identification and elimination of potential transcription factor or other undesirable sequences can be repeated until a nucleotide sequence is achieved containing a maximum number of preferred codons and a minimum number of undesired sequences including transcription regulatory sequences or other undesirable sequences. Also, optionally, desired sequences, e.g., restriction enzyme recognition sites, can be introduced. After a synthetic nucleotide sequence is designed and constructed, its properties relative to the parent nucleic acid sequence can be determined by methods well known to the art. For example, the expression of the synthetic and target nucleic acids in a series of vectors in a particular cell can be compared.

Thus, generally, the method of the invention comprises identifying a target nucleic acid sequence, and a host cell of interest, for example, a plant (dicot or monocot), fungus, yeast or mammalian cell. Preferred host cells are mammalian host cells such as CHO, COS, 293, HeLa, CV-1 and NIH3T3 cells. Based on preferred codon usage in the host cell(s) and, optionally, low codon usage in the host cell(s), e.g., high usage mammalian codons and low usage *E. coli* and mammalian codons, codons to be replaced are determined. Concurrent, subsequent or prior to selecting codons to be replaced, desired and undesired sequences, such as undesired transcriptional regulatory sequences, in the target sequence are identified. These sequences, including transcriptional regulatory sequences and restriction endonuclease sites, can be identified using databases and software such as TRANSFAC[®] (Transcription Factor Database, <http://www.gene-regulation.com/>), Match[™] (<http://www.gene-regulation.com/>),

MatInspector (Genomatix, <http://www.genomatix.de>), EPD (Eukaryotic Promoter Database, <http://www.epd.isb-sib.ch/>), REBASE[®] (Restriction Enzyme Database, NEB, <http://rebase.neb.com>), TESS (Transcription Element Search System, <http://www.cbil.upenn.edu/tess/>), MAR-Wiz (Futuresoft, <http://www.futuresoft.org>), Lasergene[®] (DNASTAR, <http://www.dnastar.com>), Vector NTI[™] (Invitrogen, <http://www.invitrogen.com>), and Sequence Manipulation Suite (<http://www.bioinformatics.org/SMS/index.html>).

Links to other databases and sequence analysis software are listed at

<http://www.expasy.org/alinks.html>. After one or more sequences are identified,

- 10 the modification(s) may be introduced. Once a desired synthetic nucleotide sequence is obtained, it can be prepared by methods well known to the art (such as nucleic acid amplification reactions with overlapping primers), and its structural and functional properties compared to the target nucleic acid sequence, including, but not limited to, percent homology, presence or absence of certain
- 15 sequences, for example, restriction sites, percent of codons changed (such as an increased or decreased usage of certain codons) and/or expression rates.

As described below, the method was used to create synthetic reporter genes encoding firefly luciferases and selectable polypeptides, and synthetic sequences for vector backbones. Synthetic sequences may support greater levels

20 of expression and/or reduced aberrant expression than the corresponding native or parent sequences for the protein. The native and parent sequences may demonstrate anomalous transcription characteristics when expressed in mammalian cells, which are likely not evident in the synthetic sequences.

25 Exemplary Uses of the Synthetic Nucleotide Sequences

The synthetic genes of the invention preferably encode the same proteins as their native counterpart (or nearly so), but have improved codon usage while being largely devoid of regulatory elements in the coding (it is recognized that a small number of amino acid changes may be desired to enhance a property of the

30 native counterpart protein, e.g. to enhance luminescence of a luciferase) and noncoding regions. This increases the level of expression of the protein the synthetic gene encodes and reduces the risk of anomalous expression of the protein. For example, studies of many important events of gene regulation,

which may be mediated by weak promoters, are limited by insufficient reporter signals from inadequate expression of the reporter proteins. Also, the use of some selectable markers may be limited by the expression of that marker in an exogenous cell. Thus, synthetic selectable marker genes which have improved
5 codon usage for that cell, and have a decrease in other undesirable sequences, (e.g., transcription factor binding sites), can permit the use of those markers in cells that otherwise were undesirable as hosts for those markers.

Promoter crosstalk is another concern when a co-reporter gene is used to normalize transfection efficiencies. With the enhanced expression of synthetic
10 genes, the amount of DNA containing strong promoters can be reduced, or DNA containing weaker promoters can be employed, to drive the expression of the co-reporter. In addition, there may be a reduction in the background expression from the synthetic reporter genes of the invention. This characteristic makes synthetic reporter genes more desirable by minimizing the sporadic expression
15 from the genes and reducing the interference resulting from other regulatory pathways.

The use of reporter genes in imaging systems, which can be used for *in vivo* biological studies or drug screening, is another use for the synthetic genes of the invention. Due to their increased level of expression, the protein encoded by
20 a synthetic gene is more readily detectable by an imaging system. In fact, using a synthetic *Renilla* luciferase gene, luminescence in transfected CHO cells was detected visually without the aid of instrumentation.

In addition, the synthetic genes may be used to express fusion proteins, for example fusions with secretion leader sequences or cellular localization
25 sequences, to study transcription in difficult-to-transfect cells such as primary cells, and/or to improve the analysis of regulatory pathways and genetic elements. Other uses include, but are not limited to, the detection of rare events that require extreme sensitivity (e.g., studying RNA recoding), use with IRES, to improve the efficiency of *in vitro* translation or *in vitro* transcription-translation
30 coupled systems such as TnT (Promega Corp., Madison, WI), study of reporters optimized to different host organisms (e.g., plants, fungus, and the like), use of multiple genes as co-reporters to monitor drug toxicity, as reporter molecules in multiwell assays, and as reporter molecules in drug screening with the advantage

of minimizing possible interference of reporter signal by different signal transduction pathways and other regulatory mechanisms.

Additionally, uses for the synthetic nucleotide sequences of the invention include fluorescence activated cell sorting (FACS), fluorescent microscopy, to
5 detect and/or measure the level of gene expression *in vitro* and *in vivo*, (e.g., to determine promoter strength), subcellular localization or targeting (fusion protein), as a marker, in calibration, in a kit (e.g., for dual assays), for *in vivo* imaging, to analyze regulatory pathways and genetic elements, and in multi-well formats.

10 Further, although reporter genes are widely used to measure transcription events, their utility can be limited by the fidelity and efficiency of reporter expression. For example, in U.S. Patent No. 5,670,356, a firefly luciferase gene (referred to as luc+) was modified to improve the level of luciferase expression. While a higher level of expression was observed, it was not determined that
15 higher expression had improved regulatory control.

The invention will be further described by the following nonlimiting examples. In particular, the synthetic nucleic acid molecules of the invention may be derived by other methods as well as by variations on the methods described herein.

20

Example 1

Synthetic Click Beetle (RD and GR) Luciferase Nucleic Acid Molecules

LucPpIYG is a wild-type click beetle luciferase that emits yellow-green luminescence (Wood, 1989). A mutant of LucPpIYG named YG#81-6G01 was
25 envisioned. YG#81-6G01 lacks a peroxisome targeting signal, has a lower K_M for luciferin and ATP, has increased signal stability and increased temperature stability when compared to the wild type (PCT/WO9914336). YG #81-6G01 was mutated to emit green luminescence by changing Ala at position 224 to Val (A224V is a green-shifting mutation), or to emit red luminescence by
30 simultaneously introducing the amino acid substitutions A224H, S247H, N346I, and H348Q (red-shifting mutation set) (PCT/WO9518853)

Using YG #81-6G01 as a parent gene, two synthetic gene sequences were designed. One codes for a luciferase emitting green luminescence (GR) and one

for a luciferase emitting red luminescence (RD). Both genes were designed to 1)
have optimized codon usage for expression in mammalian cells, 2) have a
reduced number of transcriptional regulatory sites including mammalian
transcription factor binding sites, splice sites, poly(A) sites and promoters, as
5 well as prokaryotic (*E. coli*) regulatory sites, 3) be devoid of unwanted restriction
sites, e.g., those which are likely to interfere with standard cloning procedures,
and 4) have a low DNA sequence identity compared to each other in order to
minimize genetic rearrangements when both are present inside the same cell. In
addition, desired sequences, e.g., a Kozak sequence or restriction enzyme
10 recognition sites, may be identified and introduced.

Not all design criteria could be met equally well at the same time. The
following priority was established for reduction of transcriptional regulatory
sites: elimination of transcription factor (TF) binding sites received the highest
priority, followed by elimination of splice sites and poly(A) sites, and finally
15 prokaryotic regulatory sites. When removing regulatory sites, the strategy was to
work from the lesser important to the most important to ensure that the most
important changes were made last. Then the sequence was rechecked for the
appearance of new lower priority sites and additional changes made as needed.
Thus, the process for designing the synthetic GR and RD gene sequences, using
20 computer programs described herein, involved 5 optionally iterative steps that
are detailed below

1. Optimized codon usage and changed A224V to create GRver1,
separately changed A224H, S247H, H348Q and N346I to create
RDver1. These particular amino acid changes were maintained
25 throughout all subsequent manipulations to the sequence.
2. Removed undesired restriction sites, prokaryotic regulatory sites,
splice sites, poly(A) sites thereby creating GRver2 and RDver2.
3. Removed transcription factor binding sites (first pass) and removed
any newly created undesired sites as listed in step 2 above thereby
30 creating GRver3 and RDver3.
4. Removed transcription factor binding sites created by step 3 above
(second pass) and removed any newly created undesired sites as listed
in step 2 above thereby creating GRver4 and RDver4.

5. Removed transcription factor binding sites created by step 4 above (third Pass) and confirmed absence of sites listed in step 2 above thereby creating GRver5 and RDver5.
6. Constructed the actual genes by PCR using synthetic oligonucleotides corresponding to fragments of GRver5 and RDver5 designed sequences thereby creating GR6 and RD7. GR6, upon sequencing was found to have the serine residue at amino acid position 49 mutated to an asparagine and the proline at amino acid position 230 mutated to a serine (S49N, P230S). RD7, upon sequencing was found to have the histidine at amino acid position 36 mutated to a tyrosine (H36Y). These changes occurred during the PCR process.
4. The mutations described in step 6 above (S49N, P230S for GR6 and H36Y for RD7) were reversed to create GRver5.1 and RDver5.1.
5. RDver5.1 was further modified by changing the arginine codon at position 351 to a glycine codon (R351G) thereby creating RDver5.2 with improved spectral properties compared to RDver5.1.
6. RDver5.2 was further mutated to increase luminescence intensity thereby creating RD156-1H9 which encodes four additional amino acid changes (M2I, S349T, K488T, E538V) and three silent single base changes (see U.S. application Serial No. 09/645,706, filed August 24, 2000, the disclosure of which is incorporated by reference herein).

1. Optimize codon usage and introduce mutations determining luminescence

color

The starting gene sequence for this design step was YG #81-6G01.

a) Optimize codon usage:

The strategy was to adapt the codon usage for optimal expression in human cells and at the same time to avoid *E. coli* low-usage codons. Based on these requirements, the best two codons for expression in human cells for all amino acids with more than two codons were selected (see Wada et al., 1990). In the selection of codon pairs for amino acids with six codons, the selection was biased towards pairs that have the largest number of mismatched bases to allow

design of GR and RD genes with minimum sequence identity (codon distinction):

| | | | |
|---|--------------|--------------|--------------|
| | Arg: CGC/CGT | Leu: CTG/TTG | Ser: TCT/AGC |
| | Thr: ACC/ACT | Pro: CCA/CCT | Ala: GCC/GCT |
| 5 | Gly: GGC/GGT | Val: GTC/GTG | Ile: ATC/ATT |

Based on this selection of codons, two gene sequences encoding the YG#81-6G01 luciferase protein sequence were computer generated. The two genes were designed to have minimum DNA sequence identity and at the same time closely similar codon usage. To achieve this, each codon in the two genes was replaced
 10 by a codon from the limited list described above in an alternating fashion (e.g., Arg_(n) is CGC in gene 1 and CGT in gene 2, Arg_(n+1) is CGT in gene 1 and CGC in gene 2).

For subsequent steps in the design process it was anticipated that changes had to be made to this limited optimal codon selection in order to meet other
 15 design criteria, however, the following low-usage codons in mammalian cells were not used unless needed to meet criteria of higher priority:

| | | |
|----------|----------|----------|
| Arg: CGA | Leu: CTA | Ser: TCG |
| Pro: CCG | Val: GTA | Ile: ATA |

Also, the following low-usage codons in *E. coli* were avoided when reasonable
 20 (note that 3 of these match the low-usage list for mammalian cells):

| | | |
|----------------------|----------|----------|
| Arg: CGA/CGG/AGA/AGG | | |
| Leu: CTA | Pro: CCC | Ile: ATA |

b) Introduce mutations determining luminescence color:

Into one of the two codon-optimized gene sequences was introduced the
 25 single green-shifting mutation and into the other were introduced the 4 red-shifting mutations as described above.

The two output sequences from this first design step were named GRver1 (version 1 GR) and RDver1 (version 1 RD). Their DNA sequences are 63% identical (594 mismatches), while the proteins they encode differ only by the 4
 30 amino acids that determine luminescence color (see Figures 2 and 3 for an alignment of the DNA and protein sequences).

Tables 1 and 2 show, as an example, the codon usage for valine and leucine in human genes, the parent gene YG#81-6G01, the codon-optimized

synthetic genes GRver1 and RDver1, as well as the final versions of the synthetic genes after completion of step 5 in the design process (GRver5 and RDver5).

Table 1: Valine

| Codon | Human | Parent | GR ver1 | RD ver1 | GR ver5 | RD ver5 |
|-------|-------|--------|---------|---------|---------|---------|
| GTA | 4 | 13 | 0 | 0 | 1 | 1 |
| GTC | 13 | 4 | 25 | 24 | 21 | 26 |
| GTG | 24 | 12 | 25 | 25 | 25 | 17 |
| GTT | 9 | 20 | 0 | 0 | 3 | 5 |

5

Table 2: Leucine

| Codon | Human | Parent | GR ver1 | RD ver1 | GR ver5 | RD ver5 |
|-------|-------|--------|---------|---------|---------|---------|
| CTA | 3 | 5 | 0 | 0 | 0 | 0 |
| CTC | 12 | 4 | 0 | 1 | 12 | 11 |
| CTG | 24 | 4 | 28 | 27 | 19 | 18 |
| CTT | 6 | 12 | 0 | 0 | 1 | 1 |
| TTA | 3 | 17 | 0 | 0 | 0 | 0 |
| TTG | 6 | 13 | 27 | 27 | 23 | 25 |

2. Remove undesired restriction sites, prokaryotic regulatory sites, splice sites and poly(A) sites

10 The starting gene sequences for this design step were GRver1 and RDver1.

a) Remove undesired restriction sites:

To check for the presence and location of undesired restriction sites, the sequences of both synthetic genes were compared against a database of restriction enzyme recognition sequences (REBASE ver.712,

15 <http://www.neb.com/rebase>) using standard sequence analysis software (GenePro ver 6.10, Riverside Scientific Ent.).

Specifically, the following restriction enzymes were classified as undesired:

- *Bam*H I, *Xho* I, *Sfi* I, *Kpn* I, *Sac* I, *Mlu* I, *Nhe* I, *Sma* I, *Xho* I, *Bgl* II, *Hind* III, *Nco* I, *Nar* I, *Xba* I, *Hpa* I, *Sal* I,
- 20 - other cloning sites commonly used: *Eco*R I, *Eco*R V, *Cla* I,
- eight-base cutters (commonly used for complex constructs),
- *Bst*E II (to allow N-terminal fusions),
- *Xcm* I (can generate A/T overhang used for T-vector cloning).

To eliminate undesired restriction sites when found in a synthetic gene, one or more codons of the synthetic gene sequence were altered in accordance with the codon optimization guidelines described in 1a above.

b) Remove prokaryotic (*E. coli*) regulatory sequences:

5 To check for the presence and location of prokaryotic regulatory sequences, the sequences of both synthetic genes were searched for the presence of the following consensus sequences using standard sequence analysis software (GenePro):

- TATAAT (-10 Pribnow box of promoter)
- 10 - AGGA or GGAG (ribosome binding site; only considered if paired with a methionine codon 12 or fewer bases downstream).

To eliminate such regulatory sequences when found in a synthetic gene, one or more codons of the synthetic gene at sequence were altered in accordance with the codon optimization guidelines described in 1a above.

15 **c) Remove splice sites:**

 To check for the presence and location of splice sites, the DNA strand corresponding to the primary RNA transcript of each synthetic gene was searched for the presence of the following consensus sequences (see Watson et al., 1983) using standard sequence analysis software (GenePro):

- 20 - splice donor site: AG | GTRAGT (exon | intron), the search was performed for AGGTRAG and the lower stringency GGTRAGT;
- splice acceptor site: (Y)_nNCAG | G (intron | exon), the search was performed with n = 1.

 To eliminate splice sites found in a synthetic gene, one or more codons of the
 25 synthetic gene sequence were altered in accordance with the codon optimization guidelines described in 1a above. Splice acceptor sites were generally difficult to eliminate in one gene without introducing them into the other gene because they tended to contain one of the two only Gln codons (CAG); they were removed by placing the Gln codon CAA in both genes at the expense of a slightly increased
 30 sequence identity between the two genes.

d) Remove poly(A) sites:

To check for the presence and location of poly(A) sites, the sequences of both synthetic genes were searched for the presence of the following consensus sequence using standard sequence analysis software (GenePro):

5 - AATAAA.

To eliminate each poly(A) addition site found in a synthetic gene, one or more codons of the synthetic gene sequence were altered in accordance with the codon optimization guidelines described in 1a above. The two output sequences from this second design step were named GRver2 and RDver2. Their DNA sequences
10 are 63% identical (590 mismatches).

3. Remove transcription factor (TF) binding sites, then repeat steps 2 a-d

The starting gene sequences for this design step were GRver2 and RDver2.

15 To check for the presence, location and identity of potential TF binding sites, the sequences of both synthetic genes were used as query sequences to search a database of transcription factor binding sites (TRANSFAC v3.2). The TRANSFAC database (<http://transfac.gbf.de/TRANSFAC/index.html>) holds information on gene regulatory DNA sequences (TF binding sites) and proteins
20 (TFs) that bind to and act through them. The SITE table of TRANSFAC Release 3.2 contains 4,401 entries of individual (putative) TF binding sites (including TF binding sites in eukaryotic genes, in artificial sequences resulting from mutagenesis studies and *in vitro* selection procedures based on random
25 oligonucleotide mixtures or specific theoretical considerations, and consensus binding sequences (from Faisst and Meyer, 1992).

The software tool used to locate and display these TF binding sites in the synthetic gene sequences was TESS (Transcription Element Search Software, <http://agave.humgen.upenn.edu/tess/index.html>). The filtered string-based search option was used with the following user-defined search parameters:

- 30 - Factor Selection Attribute: Organism Classification
 - Search Pattern: Mammalia
 - Max. Allowable Mismatch %: 0
 - Min. element length: 5

- Min. log-likelihood: 10

This parameter selection specifies that only mammalian TF binding sites (approximately 1,400 of the 4,401 entries in the database) that are at least 5 bases long will be included in the search. It further specifies that only TF binding sites that have a perfect match in the query sequence and a minimum log likelihood (LLH) score of 10 will be reported. The LLH scoring method assigns 2 to an unambiguous match, 1 to a partially ambiguous match (e.g., A or T match W) and 0 to a match against 'N'. For example, a search with parameters specified above would result in a "hit" (positive result or match) for TATAA (SEQ ID NO:50) (LLH = 10), STRATG (SEQ ID NO:51) (LLH = 10), and MTTNCNNMA (SEQ ID NO:52) (LLH = 10) but not for TRATG (SEQ ID NO: 53) (LLH = 9) if these four TF binding sites were present in the query sequence. A lower stringency test was performed at the end of the design process to re-evaluate the search parameters.

When TESS was tested with a mock query sequence containing known TF binding sites it was found that the program was unable to report matches to sites ending with the 3' end of the query sequence. Thus, an extra nucleotide was added to the 3' end of all query sequences to eliminate this problem.

The first search for TF binding sites using the parameters described above found about 100 transcription factor binding sites (hits) for each of the two synthetic genes (GRver2 and RDver2). All sites were eliminated by changing one or more codons of the synthetic gene sequences in accordance with the codon optimization guidelines described in 1a above. However, it was expected that some these changes created new TF binding sites, other regulatory sites, and new restriction sites. Thus, steps 2 a-d were repeated as described, and 4 new restriction sites and 2 new splice sites were removed. The two output sequences from this third design step were named GRver3 and RDver3. Their DNA sequences are 66% identical (541 mismatches).

4. Remove new transcription factor (TF) binding sites, then repeat steps 2 a-d

The starting gene sequences for this design step were GRver3 and RDver3.

This fourth step is an iteration of the process described in step 3. The search for newly introduced TF binding sites yielded about 50 hits for each of the two synthetic genes. All sites were eliminated by changing one or more codons of the synthetic gene sequences in general accordance with the codon optimization guidelines described in 1a above. However, more high to medium usage codons were used to allow elimination of all TF binding sites. The lowest priority was placed on maintaining low sequence identity between the GR and RD genes. Then steps 2 a-d were repeated as described. The two output sequences from this fourth design step were named GRver4 and RDver4. Their DNA sequences are 68% identical (506 mismatches).

5. Remove new transcription factor (TF) binding sites, then repeat steps 2 a-d

The starting gene sequences for this design step were GRver4 and RDver4.

This fifth step is another iteration of the process described in step 3 above. The search for new TF binding sites introduced in step 4 yielded about 20 hits for each of the two synthetic genes. All sites were eliminated by changing one or more codons of the synthetic gene sequences in general accordance with the codon optimization guidelines described in 1a above. However, more high to medium usage codons were used (these are all considered "preferred") to allow elimination of all TF binding sites. The lowest priority was placed on maintaining low sequence identity between the GR and RD genes. Then steps 2 a-d were repeated as described. Only one acceptor splice site could not be eliminated. As a final step the absence of all TF binding sites in both genes as specified in step 3 was confirmed. The two output sequences from this fifth and last design step were named GRver5 and RDver5. Their DNA sequences are 69% identical (504 mismatches).

Additional evaluation of GRver5 and RDver5

a) Use lower stringency parameters for TESS:

The search for TF binding sites was repeated as described in step 3 above, but with even less stringent user-defined parameters:

- setting LLH to 9 instead of 10 did not result in new hits;

- setting LLH to 0 through 8 (incl.) resulted in hits for two additional sites, MAMAG (22 hits) and CTKTK (24 hits);
- setting LLH to 8 and the minimum element length to 4, the search yielded (in addition to the two sites above) different 4-base sites for AP-1, NF-1, and c-Myb that are shortened versions of their longer
5 respective consensus sites which were eliminated in steps 3-5 above.

It was not realistic to attempt complete elimination of these sites without introduction of new sites, so no further changes were made.

b) Search different database:

- 10 The Eukaryotic Promoter Database (release 45) contains information about reliably mapped transcription start sites (1253 sequences) of eukaryotic genes. This database was searched using BLASTN 1.4.11 with default parameters (optimized to find nearly identical sequences rapidly; see Altschul et al, 1990) at the National Center for Biotechnology Information site
15 (<http://www.ncbi.nlm.nih.gov/cgi-bin/BLAST>). To test this approach, a portion of pGL3-Control vector sequence containing the SV40 promoter and enhancer was used as a query sequence, yielding the expected hits to SV40 sequences. No hits were found when using the two synthetic genes as query sequences.

20 Summary of GRver5 and RDver5 synthetic gene properties

Both genes, which at this stage were still only "virtual" sequences in the computer, have a codon usage that strongly favors mammalian high-usage codons and minimizes mammalian and *E. coli* low-usage codons.

- Both genes are also completely devoid of eukaryotic TF binding sites
25 consisting of more than four unambiguous bases, donor and acceptor splice sites (one exception: GRver5 contains one splice acceptor site), poly(A) sites, specific prokaryotic (*E. coli*) regulatory sequences, and undesired restriction sites.

- The gene sequence identity between GRver5 and RDver5 is only 69% (504 base mismatches) while their encoded proteins are 99% identical (4 amino
30 acid mismatches). Their identity with the parent sequence YG#81-6G1 is 74% (GRver5) and 73% (RDver5). Their base composition is 49.9% GC (GRver5) and 49.5% GC (RDver5), compared to 40.2% GC for the parent YG#81-6G01.

Construction of synthetic genes

The two synthetic genes were constructed by assembly from synthetic oligonucleotides in a thermocycler followed by PCR amplification of the full-length genes (similar to Stemmer et al. (1995) Gene. 164, pp. 49-53).

- 5 Unintended mutations that interfered with the design goals of the synthetic genes were corrected.

a) Design of synthetic oligonucleotides:

- 10 The synthetic oligonucleotides were mostly 40mers that collectively code for both complete strands of each designed gene (1,626 bp) plus flanking regions needed for cloning (1,950 bp total for each gene). The 5' and 3' boundaries of all oligonucleotides specifying one strand were generally placed in a manner to give an average offset/overlap of 20 bases relative to the boundaries of the oligonucleotides specifying the opposite strand.

- 15 The ends of the flanking regions of both genes matched the ends of the amplification primers (pRAMtailup: 5'-gtactgagacgacgccagcccaagccttaggcctgagtg SEQ ID NO:54, and pRAMtaildn: 5'-ggcatgagcgtgaactgactgaactagcggccgcccag SEQ ID NO:55) to allow cloning of the genes into our *E. coli* expression vector pRAM (WO99/14336).

- 20 A total of 183 oligonucleotides were designed: fifteen oligonucleotides that collectively encode the upstream and downstream flanking sequences and 168 oligonucleotides (4 x 42) that encode both strands of the two genes.

- 25 All 183 oligonucleotides were run through the hairpin analysis of the OLIGO software (OLIGO 4.0 Primer Analysis Software © 1989-1991 by Wojciech Rychlik) to identify potentially detrimental intra-molecular loop formation. The guidelines for evaluating the analysis results were set according to recommendations of Dr. Sims (Sigma-Genosys Custom Gene Synthesis Department): oligos forming hairpins with $\Delta G < -10$ have to be avoided, those forming hairpins with $\Delta G \leq -7$ involving the 3' end of the oligonucleotide should also be avoided, while those with an overall $\Delta G \leq -5$ should not pose a problem
- 30 for this application. The analysis identified 23 oligonucleotides able to form hairpins with a ΔG between -7.1 and -4.9. Of these, 5 had blocked or nearly blocked 3' ends (0-3 free bases) and were re-designed by removing 1-4 bases at their 3' end and adding it to the adjacent oligonucleotide.

The 40mer oligonucleotide covering the sequence complementary to the poly(A) tail had a very low complexity 3' end (13 consecutive T bases). An additional 40mer was designed with a high complexity 3' end but a consequently reduced overlap with one of its complementary oligonucleotides (11 instead of 20 bases) on the opposite strand.

Even though the oligonucleotides were designed for use in a thermocycler-based assembly reaction, they could also be used in a ligation-based protocol for gene construction. In this approach, the oligonucleotides are annealed in a pairwise fashion and the resulting short double-stranded fragments are ligated using the sticky overhangs. However, this would require that all oligonucleotides be phosphorylated.

b) Gene assembly and amplification

In a first step, each of the two synthetic genes was assembled in a separate reaction from 98 oligonucleotides. The total volume for each reaction was 50 μ l:

0.5 μ M oligonucleotides (= 0.25 pmoles of each oligo)

1.0 U *Taq* DNA polymerase

0.02 U *Pfu* DNA polymerase

2 mM $MgCl_2$

0.2 mM dNTPs (each)

0.1% gelatin

Cycling conditions: (94°C for 30 seconds, 52°C for 30 seconds, and 72°C for 30 seconds) x 55 cycles.

In a second step, each assembled synthetic gene was amplified in a separate reaction. The total volume for each reaction was 50 μ l:

2.5 l assembly reaction

5.0 U *Taq* DNA polymerase

0.1 U *Pfu* DNA polymerase

1 M each primer (pRAMtailup, pRAMtaildn)

2 mM $MgCl_2$

0.2 mM dNTPs (each)

Cycling conditions: (94°C for 20 seconds, 65°C for 60 seconds, 72°C for 3 minutes) x 30 cycles.

The assembled and amplified genes were subcloned into the pRAM vector and expressed in *E. coli*, yielding 1-2% luminescent GR or RD clones. Five GR and five RD clones were isolated and analyzed further. Of the five GR clones, three had the correct insert size, of which one was weakly luminescent and one had an altered restriction pattern. Of the five RD clones, two had the correct size insert with an altered restriction pattern and one of those was weakly luminescent. Overall, the analysis indicated the presence of a large number of mutations in the genes, most likely the result of errors introduced in the assembly and amplification reactions.

10 c) Corrective assembly and amplification

To remove the large number of mutations present in the full-length synthetic genes we performed an additional assembly and amplification reaction for each gene using the proof-reading DNA polymerase *Tli*. The assembly reaction contained, in addition to the 98 GR or RD oligonucleotides, a small amount of DNA from the corresponding full-length clones with mutations described above. This allows the oligos to correct mutations present in the templates.

The following assembly reaction was performed for each of the synthetic genes. The total volume for each reaction was 50 μ l:

| | |
|----|---|
| 20 | 0.5 μ M oligonucleotides (= 0.25 pmoles of each oligo) |
| | 0.016 pmol plasmid (mix of clones with correct insert size) |
| | 2.5 U <i>Tli</i> DNA polymerase |
| | 2 mM $MgCl_2$ |
| 25 | 0.2 mM dNTPs (each) |
| | 0.1 % gelatin |
| | Cycling conditions: 94°C for 30 seconds, then (94°C for 30 seconds, 52°C for 30 seconds, 72°C for 30 seconds) for 55 cycles, then 72°C for 5 minutes. |

30 The following amplification reaction was performed on each of the assembly reactions. The total volume for each amplification reaction was 50 μ l:

1-5 μ l of assembly reaction
40 pmol each primer (pRAMtailup, pRAMtaildn)

2.5 U *Tli* DNA polymerase

2 mM MgCl₂

0.2 mM dNTPs (each)

5 Cycling conditions: 94°C for 30 seconds, then (94°C for
20 seconds, 65°C for 60 seconds and 72°C for 3 minutes)
for 30 cycles, then 72°C for 5 minutes.

10 The genes obtained from the corrective assembly and amplification step
were subcloned into the pRAM vector and expressed in *E. coli*, yielding 75%
luminescent GR or RD clones. Forty-four GR and 44 RD clones were analyzed
with the screening robot described in WO99/14336. The six best GR and RD
clones were manually analyzed and one best GR and RD clone was selected
(GR6 and RD7). Sequence analysis of GR6 revealed two point mutations in the
coding region, both of which resulted in an amino acid substitution (S49N and
P230S). Sequence analysis of RD7 revealed three point mutations in the coding
15 region, one of which resulted in an amino acid substitution (H36Y). It was
confirmed that none of the silent point mutations introduced any regulatory or
restriction sites conflicting with the overall design criteria for the synthetic
genes.

d) Reversal of unintended amino acid substitutions

20 The unintended amino acid substitutions present in the GR6 and RD7
synthetic genes were reversed by site-directed mutagenesis to match the GRver5
and RDver5 designed sequences, thereby creating GRver5.1 and RDver5.1. The
DNA sequences of the mutated regions were confirmed by sequence analysis.

e) Improve spectral properties

25 The RDver5.1 gene was further modified to improve its spectral
properties by introducing an amino change (R351G), thereby creating RDver5.2

pGL3 vectors with RD and GR genes

30 The parent click beetle luciferase YG#81-6G1 ("YG"), and the synthetic
click beetle luciferase genes GRver5.1 ("GR"), RDver5.2 ("RD"), and RD156-
1H9 were cloned into the four pGL3 reporter vectors (Promega Corp.):

- pGL3-Basic = no promoter, no enhancer
- pGL3-Control = SV40 promoter, SV40 enhancer

- pGL3-Enhancer = SV40 enhancer (3' to luciferase coding sequences)
- pGL3-Promoter = SV40 promoter.

The primers employed in the assembly of GR and RD synthetic genes facilitated the cloning of those genes into pRAM vectors. To introduce the genes into

5 pGL3 vectors (Promega Corp., Madison, WI) for analysis in mammalian cells, each gene in a pRAM vector (pRAM RDver5.1, pRAM GRver5.1, and pRAM RD156-1H9) was amplified to introduce an *Nco* I site at the 5' end and an *Xba* I site at the 3' end of the gene. The primers for pRAM RDver5.1 and pRAM GRver5.1 were:

10 GR→5' GGA TCC CAT GGT GAA GCG TGA GAA 3' (SEQ ID NO:56) or
RD→5' GGA TCC CAT GGT GAA ACG CGA 3' (SEQ ID NO:57) and
5' CTA GCT TTT TTT TCT AGA TAA TCA TGA AGA C 3' (SEQ ID NO:58)

The primers for pRAM RD156-1H9 were:

5' GCG TAG CCA TGG TAA AGC GTG AGA AAA ATG TC 3' (SEQ ID NO:
15 59) and
5' CCG ACT CTA GAT TAC TAA CCG CCG GCC TTC ACC 3' (SEQ ID NO:
60)

The PCR included:

| | |
|----|---|
| | 100 ng DNA plasmid |
| 20 | 1 μ M primer upstream |
| | 1 μ M primer downstream |
| | 0.2 mM dNTPs |
| | 1X buffer (Promega Corp.) |
| | 5 units <i>Pfu</i> DNA polymerase (Promega Corp.) |
| 25 | Sterile nanopure H ₂ O to 50 μ l |

The cycling parameters were: 94°C for 5 minutes; (94°C for 30 seconds; 55°C for 1 minute; and 72°C for 3 minutes) x 15 cycles. The purified PCR product was digested with *Nco* I and *Xba* I, ligated with pGL3-control that was also digested with *Nco* I and *Xba* I, and the ligated products introduced to *E. coli*.

30 To insert the luciferase genes into the other pGL3 reporter vectors (basic, promoter and enhancer), the pGL3-control vectors containing each of the luciferase genes was digested with *Nco* I and *Xba* I, ligated with other pGL3

vectors that also were digested with *Nco* I and *Xba* I, and the ligated products introduced to *E. coli*. Note that the polypeptide encoded by GRver5.1 and RDver5.1 (and RD156-1H9, see below) nucleic acid sequences in pGL3 vectors has an amino acid substitution at position 2 to valine as a result of the *Nco* I site at the initiation codon in the oligonucleotide.

Because of internal *Nco* I and *Xba* I sites, the native gene in YG #81-6G01 was amplified from a *Hind* III site upstream to a *Hpa* I site downstream of the coding region and which included flanking sequences found in the GR and RD clones. The upstream primer (5'-CAA AAA GCT TGG CAT TCC GGT ACT GTT GGT AAA GCC ACC ATG GTG AAG CGA GAG- 3'; SEQ ID NO:61) and a downstream primer (5'- CAA TTG TTG TTG TTA ACT TGT TTA TT -3'; SEQ ID NO:62) were mixed with YG#81-6G01 and amplified using the PCR conditions above. The purified PCR product was digested with *Nco* I and *Xba* I, ligated with pGL3-control that was also digested with *Hind* III and *Hpa* I, and the ligated products introduced into *E. coli*. To insert YG#81-6G01 into the other pGL3 reporter vectors (basic, promoter and enhancer), the pGL3-control vectors containing YG#81-6G01 were digested with *Nco* I and *Xba* I, ligated with the other pGL3 vectors that also were digested with *Nco* I and *Xba* I, and the ligated products introduced to *E. coli*. Note that the clone of YG#81-6G01 in the pGL3 vectors has a C instead of an A at base 786, which yields a change in the amino acid sequence at residue 262 from Phe to Leu. To determine whether the altered amino acid at position 262 affected the enzyme biochemistry, the clone of YG#81-6G01 was mutated to resemble the original sequence. Both clones were then tested for expression in *E. coli*, physical stability, substrate binding, and luminescence output kinetics. No significant differences were found.

Partially purified enzymes expressed from the synthetic genes and the parent gene were employed to determine K_m for luciferin and ATP (see Table 3).

Table 3

| Enzyme | K _M (LH ₂) | K _M (ATP) |
|-----------|-----------------------------------|----------------------|
| YG parent | 2 μ M | 17 μ M |
| GR | 1.3 μ M | 25 μ M |
| RD | 24.5 μ M | 46 μ M |

In vitro eukaryotic transcription/translation reactions were also conducted using Promega's TNT T7 Quick system according to manufacturer's instructions. Luminescence levels were 1 to 37-fold and 1 to 77-fold higher (depending on the reaction time) for the synthetic GR and RD genes, respectively, compared to the parent gene (corrected for luminometer spectral sensitivity).

To test whether the synthetic click beetle luciferase genes and the wild type click beetle gene have improved expression in mammalian cells, each of the synthetic genes and the parent gene was cloned into a series of pGL3 vectors and introduced into CHO cells (Table 8). In all cases, the synthetic click beetle genes exhibited a higher expression than the native gene. Specifically, expression of the synthetic GR and RD genes was 1900-fold and 40-fold higher, respectively, than that of the parent (transfection efficiency normalized by comparison to native *Renilla* luciferase gene). Moreover, the data (basic versus control vector) show that the synthetic genes have reduced basal level transcription.

Further, in experiments with the enhancer vector where the percentage of activity in reference to the control is compared between the native and synthetic gene, the data showed that the synthetic genes have reduced risk of anomalous transcription characteristics. In particular, the parent gene appeared to contain one or more internal transcriptional regulatory sequences that are activated by the enhancer in the vector, and thus is not suitable as a reporter gene while the synthetic GR and RD genes showed a clean reporter response (transfection efficiency normalized by comparison to native *Renilla* luciferase gene). See Table 8.

Example 2

Synthetic *Renilla* Luciferase Nucleic Acid Molecule

The synthetic *Renilla* luciferase genes prepared include 1) an introduced Kozak sequence, 2) codon usage optimized for mammalian (human) expression, 3) a reduction or elimination of unwanted restriction sites, 4) removal of prokaryotic regulatory sites (ribosome binding site and TATA box), 5) removal of splice sites and poly(A) sites, and 6) a reduction or elimination of mammalian transcriptional factor binding sequences.

The process of computer-assisted design of synthetic *Renilla* luciferase genes by iterative rounds of codon optimization and removal of transcription factor binding sites and other regulatory sites as well as restriction sites can be described in three steps:

1. Using the wild type *Renilla* luciferase gene as the parent gene, codon usage was optimized, one amino acid was changed (T→A) to generate a Kozak consensus sequence, and undesired restriction sites were eliminated thereby creating synthetic gene Rlucver1.
2. Remove prokaryotic regulatory sites, splice sites, poly(A) sites and transcription factor (TF) binding sites (first pass). Then remove newly created TF binding sites. Then remove newly created undesired restriction enzyme sites, prokaryotic regulatory sites, splice sites, and poly(A) sites without introducing new TF binding sites. This thereby created Rlucver2.
3. Change 3 bases of Rlucver2 thereby creating Rluc-final.
4. The actual gene was then constructed from synthetic oligonucleotides corresponding to the Rluc-final designed sequence. All mutations resulting from the assembly or PCR process were corrected. This gene is Rluc-final.

Codon Selection

Starting with the *Renilla reniformis* luciferase sequence in Genbank (Accession No. M63501), codons were selected based on codon usage for optimal expression in human cells and to avoid *E. coli* low-usage codons. The best codon for expression in human cells (or the best two codons if found at a similar frequency) was chosen for all amino acids with more than one codon (Wada et al., 1990):

| | | |
|---|--------------|----------|
| | Arg: CGC | Lys: AAG |
| | Leu: CTG | Asn: AAC |
| | Ser: TCT/AGC | Gln: CAG |
| | Thr: ACC | His: CAC |
| 5 | Pro: CCA/CCT | Glu: GAG |
| | Ala: GCC | Asp: GAC |
| | Gly: GGC | Tyr: TAC |
| | Val: GTG | Cys: TGC |
| | Ile: ATC/ATT | Phe: TTC |

10 In cases where two codons were selected for one amino acid, they were used in an alternating fashion. To meet other criteria for the synthetic gene, the initial optimal codon selection was modified to some extent later. For example, introduction of a Kozak sequence required the use of GCT for Ala at amino acid position 2 (see below).

15 The following low-usage codons in mammalian cells were not used unless needed: Arg: CGA, CGU; Leu: CTA, UUA; Ser: TCG; Pro: CCG; Val: GTA; and Ile: ATA. The following low-usage codons in *E. coli* were also avoided when reasonable (note that 3 of these match the low-usage list for mammalian cells): Arg: CGA/CGG/AGA/AGG, Leu: CTA; Pro: CCC; Ile: 20 ATA.

Introduction of Kozak Sequences

The Kozak sequence: 5' aaccATGGCT 3' (SEQ ID NO: 63) (the *Nco* I site is underlined, the coding region is shown in capital letters) was introduced to the synthetic *Renilla* luciferase gene. The introduction of the Kozak sequence 25 changes the second amino acid from Thr to Ala (GCT).

Removal of undesired restriction sites

REBASE ver. 808 (updated August 1, 1998; Restriction Enzyme Database; www.neb.com/rebase) was employed to identify undesirable restriction sites as described in Example 1. The following undesired restriction 30 sites (in addition to those described in Example 1) were removed according to the process described in Example 1: *Eco*ICR I, *Nde*I, *Nsi*I, *Sph*I, *Spe*I, *Xma*I, *Pst*I.

The version of *Renilla* luciferase (Rluc) which incorporates all these changes is Rlucver1.

Removal of prokaryotic (*E. coli*) regulatory sequences, splice sites, and poly(A) sites

The priority and process for eliminating transcription regulation sites was as described in Example 1.

5 Removal of TF binding sites

The same process, tools, and criteria were used as described in Example 1, however, the newer version 3.3 of the TRANSFAC database was employed.

After removing prokaryotic regulatory sequences, splice sites and poly(A) sites from Rlucver1, the first search for TF binding sites identified about 60 hits.

10 All sites were eliminated with the exception of three that could not be removed without altering the amino acid sequence of the synthetic *Renilla* gene:

1. site at position 63 composed of two codons for W (TGGTGG), for CAC-binding protein T00076;
2. site at position 522 composed of codons for KMV (AAN ATG GTN), for myc-DF1 T00517;
- 15 3. site at position 885 composed of codons for EMG (GAR ATG GGN), for myc-DF1 T00517.

The subsequent second search for (newly introduced) TF binding sites yielded about 20 hits. All new sites were eliminated, leaving only the three sites described above. Finally, any newly introduced restriction sites, prokaryotic regulatory sequences, splice sites and poly(A) sites were removed without introducing new TF binding sites if possible.

Rlucver2 was obtained.

25 As in Example 1, lower stringency search parameters were specified for the TESS filtered string search to further evaluate the synthetic *Renilla* gene.

With the LLH reduced from 10 to 9 and the minimum element length reduced from 5 to 4, the TESS filtered string search did not show any new hits. When, in addition to the parameter changes listed above, the organism classification was expanded from "mammalia" to "chordata", the search yielded only four more TF binding sites. When the Min LLH was further reduced to between 8 and 0, the search showed two additional 5-base sites (MAMAG and CTKTK) which combined had four matches in Rlucver2, as well as several 4-base sites. Also as in Example 1, Rlucver2 was checked for hits to entries in the

EPD (Eukaryotic Promoter Database, Release 45). Three hits were determined one to *Mus musculus* promoter H-2L^d (*Cell*, 44, 261 (1986)), one to Herpes Simplex Virus type 1 promoter b'g'2.7 kb, and one to *Homo sapiens* DHFR promoter (*J. Mol. Biol.*, 176, 169 (1984)). However, no further changes were
 5 made to Rlucver2.

Summary of Properties for Rlucver2

- All 30 low usage codons were eliminated. The introduction of a Kozak sequence changed the second amino acid from Thr to Ala;
- 10 - base composition: 55.7% GC (*Renilla* wild-type parent gene: 36.5%);
- one undesired restriction site could not be eliminated: *EcoR* V at position 488;
- the synthetic gene had no prokaryotic promoter sequence but one potentially functional ribosome binding site (RBS) at positions 867-73
 15 (about 13 bases upstream of a Met codon) could not be eliminated;
- all poly(A) sites were eliminated;
- splice sites: 2 donor splice sites could not be eliminated (both share the amino acid sequence MGK);
- TF sites: all sites with a consensus of >4 unambiguous bases were
 20 eliminated (about 280 TF binding sites were removed) with 3 exceptions due to the preference to avoid changes to the amino acid sequence.

When introduced into pGL3, Rluc-final has a Kozak sequence (CACCATGGCT; SEQ ID NO:65). The changes in Rluc-final relative to Rlucver2 were introduced during gene assembly. One change was at position
 25 619, a C to an A, which eliminated a eukaryotic promoter sequence and reduced the stability of a hairpin structure in the corresponding oligonucleotide employed to assemble the gene. Other changes included a change from CGC to AGA at positions 218-220 (resulted in a better oligonucleotide for PCR).

30 Gene Assembly Strategy

The gene assembly protocol employed for the synthetic *Renilla* luciferase was similar to that described in Example 1.

Sense Strand primer:

5' AACCATGGCTTCCAAGGTGTACGACCCCGAGCAACGCAAA 3' (SEQ ID NO:66)

Anti-sense Strand primer:

5 5' GCTCTAGAATTACTGCTCGTTCTTCAGCACGCGCTCCACG 3' (SEQ ID NO:67)

The resulting synthetic gene fragment was cloned into a pRAM vector using *Nco* I and *Xba* I. Two clones having the correct size insert were sequenced. Four to six mutations were found in the synthetic gene from each clone. These mutations were fixed by site-directed mutagenesis (Gene Editor from Promega Corp., Madison, WI) and swapping the correct regions between these two genes. The corrected gene was confirmed by sequencing.

Other Vectors

15 To prepare an expression vector for the synthetic *Renilla* luciferase gene in a pGL-3 control vector backbone, 5 µg of pGL3-control was digested with *Nco* I and *Xba* I in 50 µl final volume with 2 µl of each enzyme and 5 µl 10X buffer B (nanopure water was used to fill the volume to 50 µl). The digestion reaction was incubated at 37°C for 2 hours, and the whole mixture was run on a 20 1% agarose gel in 1XTAE. The desired vector backbone fragment was purified using Qiagen's QIAquick gel extraction kit.

The native *Renilla* luciferase gene fragment was cloned into pGL3-control vector using two oligonucleotides, *Nco* I-RL-F and *Xba* I-RL-R, to PCR amplify native *Renilla* luciferase gene using pRL-CMV as the template. The 25 sequence for *Nco* I-RL-F is 5'-CGCTAGCCATGGCTTCGAAAGTTTATGATCC -3' (SEQ ID NO:68); the sequence for *Xba* I-RL-R is 5' GGCCAGTAACTCTAGAATTATTGTT-3' (SEQ ID NO:69). The PCR reaction was carried out as follows:

30 Reaction mixture (for 100 µl):

| | |
|------------------------|----------------------------|
| DNA template (Plasmid) | 1.0 µl (1.0 ng/µl final) |
| 10 X Rec. Buffer | 10.0 µl (Stratagene Corp.) |

- | | | |
|---|---------------------------|-------------------------------------|
| | dNTPs (25 mM each) | 1.0 µl (final 250 µM) |
| | Primer 1 (10 µM) | 2.0 µl (0.2 µM final) |
| 5 | Primer 2 (10 µM) | 2.0 µl (0.2 µM final) |
| | <i>Pfu</i> DNA Polymerase | 2.0 µl (2.5 U/µl, Stratagene Corp.) |
| | | 82.0 µl double distilled water |
- 10 PCR Reaction: heat 94°C for 2 minutes; (94°C for 20 seconds; 65°C for 1 minute; 72°C for 2 minutes; then 72°C for 5 minutes) x 25 cycles, then incubate on ice. The PCR amplified fragment was cut from a gel, and the DNA purified and stored at -20°C.
- 15 To introduce native *Renilla* luciferase gene fragment into pGL3-control vector, 5 µg of the PCR product of the native *Renilla* luciferase gene (RAM-RL-synthetic) was digested with *Nco* I and *Xba* I. The desired *Renilla* luciferase gene fragment was purified and stored at -20°C.
- Then 100 ng of insert and 100 ng of pGL3-control vector backbone were
- 20 digested with restriction enzymes *Nco* I and *Xba* I and ligated together. Then 2 µl of the ligation mixture was transformed into JM109 competent cells. Eight ampicillin resistance clones were picked and their DNA isolated. DNA from each positive clone of pGL3-control-native and pGL3-control-synthetic was purified. The correct sequences for the native gene and the synthetic gene in the
- 25 vectors were confirmed by DNA sequencing.
- To determine whether the synthetic *Renilla* luciferase gene has improved expression in mammalian cells, the gene was cloned into the mammalian expression vector pGL3-control vector under the control of SV40 promoter and SV40 early enhancer. The native *Renilla* luciferase gene was also cloned into
- 30 the pGL-3 control vector so that the expression from synthetic gene and the native gene could be compared. The expression vectors were then transfected into four common mammalian cell lines (CHO, NIH3T3, Hela and CV-1; Table 9), and the expression levels compared between the vectors with the synthetic gene versus the native gene. The amount of DNA used was at two different
- 35 levels to ascertain that expression from the synthetic gene is consistently

increased at different expression levels. The results show a 70-600 fold increase of expression for the synthetic *ReniZla* luciferase gene in these cells (Table 4).

Table 4

| <u>Cell Type</u> | <u>Amount Vector</u> | <u>Fold Expression Increase</u> |
|------------------|----------------------|---------------------------------|
| CHO | 0.2 µg | 142 |
| | 2.8 µg | 145 |
| NIH3T3 | 0.2 µg | 326 |
| | 2.0 µg | 593 |
| HeLa | 0.2 µg | 185 |
| | 1.0 µg | 103 |
| CV-1 | 0.2 µg | 68 |
| | 2.0 µg | 72 |

5

One important advantage of luciferase reporter is its short protein half-life. The enhanced expression could also result from extended protein half-life and, if so, this gives an undesired disadvantage of the new gene. This possibility is ruled out by a cycloheximide chase ("CHX Chase") experiment, which demonstrated that there was no increase of protein half-life resulted from the humanized *Renilla* luciferase gene.

To ensure that the increase in expression is not limited to one expression vector backbone, is promoter specific and/or cell specific, a synthetic *Renilla* gene (Rluc-final) as well as native *Renilla* gene were cloned into different vector backbones and under different promoters. The synthetic gene always exhibited increased expression compared to its wild-type counterpart (Table 5).

15

Table 5

| <u>Vector</u> | NIH-3T3 | HeLa | CHO |
|--------------------|-----------|-----------|-----------|
| pRL-tk, native | 3,834.6 | 922.4 | 7,671.9 |
| pRL-tk, synthetic | 13,252.5 | 9,040.2 | 41,743.5 |
| pRL-CMV, native | 168,062.2 | 842,482.5 | 153,539.5 |
| pRL-CMV, synthetic | 2,168,129 | 8,440,306 | 2,532,576 |
| pRL-SV40, native | 224,224.4 | 346,787.6 | 85,323.6 |

| <u>Vector</u> | NIH-3T3 | HeLa | CHO |
|-----------------------------|-----------|-----------|-----------|
| pRL-SV40, synthetic | 1,469,588 | 2,632,510 | 1,422,830 |
| pRL-null, native | 2,853.8 | 431.7 | 2,434 |
| pRL-null, synthetic | 9,151.17 | 2,439 | 28,317.1 |
| pRGL3b, native | 12 | 21.8 | 17 |
| pRGL3b, synthetic | 130.5 | 212.4 | 1,094.5 |
| pRGL3-tk, native | 27.9 | 155.5 | 186.4 |
| pRGL3-tk, synthetic | 6,778.2 | 8,782.5 | 9,685.9 |
| pRL-tk no intron, native | 31.8 | 165 | 93.4 |
| pRL-tk no intron, synthetic | 6,665.5 | 6,379 | 21,433.1 |

Table 6Percent of control vector

| <u>Vector</u> | <u>CHO cells</u> | <u>NIH3T3 cells</u> | <u>HeLa cells</u> |
|------------------------|------------------|---------------------|-------------------|
| pRL-control native | 100 | 100 | 100 |
| pRL-control synthetic | 100 | 100 | 100 |
| pRL-basic native | 4.1 | 5.6 | 0.2 |
| pRL-basic synthetic | 0.4 | 0.1 | 0.0 |
| pRL-promoter native | 5.9 | 7.8 | 0.6 |
| pRL-promoter synthetic | 15.0 | 9.9 | 1.1 |
| pRL-enhancer native | 42.1 | 123.9 | 52.7 |
| pRL-enhancer synthetic | 2.6 | 1.5 | 5.4 |

With reduced spurious expression the synthetic gene should exhibit less basal level transcription in a promoterless vector. The synthetic and native *Renilla* luciferase genes were cloned into the pGL3-basic vector to compare the basal level of transcription. Because the synthetic gene itself has increased expression efficiency, the activity from the promoterless vector cannot be compared directly to judge the difference in basal transcription, rather, this is taken into consideration by comparing the percentage of activity from the

promoterless vector in reference to the control vector (expression from the basic vector divided by the expression in the fully functional expression vector with both promoter and enhancer elements). The data demonstrate that the synthetic *Renilla* luciferase has a lower level of basal transcription than the native gene in mammalian cells (Table 6).

It is well known to those skilled in the art that an enhancer can substantially stimulate promoter activity. To test whether the synthetic gene has reduced risk of inappropriate transcriptional characteristics, the native and synthetic gene were introduced into a vector with an enhancer element (pGL3-enhancer vector). Because the synthetic gene has higher expression efficiency, the activity of both cannot be compared directly to compare the level of transcription in the presence of the enhancer, however, this is taken into account by using the percentage of activity from enhancer vector in reference to the control vector (expression in the presence of enhancer divided by the expression in the fully functional expression vector with both promoter and enhancer elements). Such results show that when native gene is present, the enhancer alone is able to stimulate transcription from 42-124% of the control, however, when the native gene is replaced by the synthetic gene in the same vector, the activity only constitutes 1-5% of the value when the same enhancer and a strong SV40 promoter are employed. This clearly demonstrates that synthetic gene has reduced risk of spurious expression (Table 6).

The synthetic *Renilla* gene (Rluc-final) was used in *in vitro* systems to compare translation efficiency with the native gene. In a T7 quick coupled transcription/translation system (Promega Corp., Madison, WI), pRL-null native plasmid (having the native *Renilla* luciferase gene under the control of the T7 promoter) or the same amount of pRL-null-synthetic plasmid (having the synthetic *Renilla* luciferase gene under the control of the T7 promoter) was added to the TNT reaction mixture and luciferase activity measured every 5 minutes up to 60 minutes. Dual Luciferase assay kit (Promega Corp.) was used to measure *Renilla* luciferase activity. The data showed that improved expression was obtained from the synthetic gene. To further evidence the increased translation efficiency of the synthetic gene, RNA was prepared by an *in vitro* transcription system, then purified. pRL-null (native or synthetic) vectors

were linearized with *Bam*HI. The DNA was purified by multiple phenol-chloroform extraction followed by ethanol precipitation. An *in vitro* T7 transcription system was employed by prepare RNAs. The DNA template was removed by using RNase-free DNase, and RNA was purified by phenol-chloroform extraction followed by multiple isopropanol precipitations. The same amount of purified RNA, either for the synthetic gene or the native gene, was then added to a rabbit reticulocyte lysate or wheat germ lysate. Again, the synthetic *Renilla* luciferase gene RNA produced more luciferase than the native one. These data suggest that the translation efficiency is improved by the synthetic sequence. To determine why the synthetic gene was highly expressed in wheat germ, plant codon usage was determined. The lowest usage codons in higher plants coincided with those in mammals.

Reporter gene assays are widely used to study transcriptional regulation events. This is often carried out in co-transfection experiments, in which, along with the primary reporter construct containing the testing promoter, a second control reporter under a constitutive promoter is transfected into cells as an internal control to normalize experimental variations including transfection efficiencies between the samples. Control reporter signal, potential promoter cross talk between the control reporter and primary reporter, as well as potential regulation of the control reporter by experimental conditions, are important aspects to consider for selecting a reliable co-reporter vector.

As described above, vector constructs were made by cloning synthetic *Renilla* luciferase gene into different vector backbones under different promoters. All the constructs showed higher expression in the three mammalian cell lines tested (Table 5). Thus, with better expression efficiency, the synthetic *Renilla* luciferase gives out higher signal when transfected into mammalian cells.

Because a higher signal is obtained, less promoter activity is required to achieve the same reporter signal, this reduced risk of promoter interference. CHO cells were transfected with 50 ng pGL3-control (firefly *luc*+) plus one of 5 different amounts of native pRL-TK plasmid (50, 100, 500, 1000, or 2000 ng) or synthetic pRL-TK (5, 10, 50, 100, or 200 ng). To each transfection, pUC19 carrier DNA was added to a total of 3 μ g DNA. 10 fold less pRL-TK DNA gave

similar or more signal as the native gene, with reduced risk of inhibiting expression from the primary reporter pGL3-control.

Experimental treatment sometimes may activate cryptic sites within the gene and cause induction or suppression of the co-reporter expression, which would compromise its function as co-reporter for normalization of transfection efficiencies. One example is that TPA induces expression of co-reporter vectors harboring the wild-type gene when transfecting MCF-7 cells. 500 ng pRL-TK (native), 5 µg native and synthetic pRG-B, 2.5 µg native and synthetic pRG-TK were transfected per well of MCF-7 cells. 100 ng/well pGL3-control (firefly luc+) was co-transfected with all RL plasmids. Carrier DNA, pUC19, was used to bring the total DNA transfected to 5.1 µg/well. 15.3 µl TransFast Transfection Reagent (Promega Corp., Madison, WI) was added per well. Sixteen hours later, cells were trypsinized, pooled and split into six wells of a 6-well dish and allowed to attach to the well for 8 hours. Three wells were then treated with the 0.2 nM of the tumor promoter, TPA (phorbol-12-myristate-13-acetate, Calbiochem #524400-S), and three wells were mock treated with 20 µl DMSO. Cells were harvested with 0.4 ml Passive Lysis Buffer 24 hours post TPA addition. The results showed that by using the synthetic gene, undesirable change of co-reporter expression by experimental stimuli can be avoided (Table 7). This demonstrates that using synthetic gene can reduce the risk of anomalous expression.

Table 7

| <u>Vector</u> | Rlu | Fold Induction |
|-----------------------------|-----|----------------|
| pRL-tk untreated (native) | 184 | |
| pRL-tk TPA treated (native) | 812 | 4.4 |
| pRG-B untreated (native) | 1 | |
| pRG-B TPA treated (native) | 8 | 8.0 |
| pRG-B untreated (final) | 132 | |
| pRG-B TPA treated (final) | 195 | 1.47 |
| pRG-tk untreated (native) | 44 | |

| <u>Vector</u> | <u>Rlu</u> | <u>Fold Induction</u> |
|-----------------------------|------------|-----------------------|
| pRG-tk TPA treated (native) | 192 | 4.36 |
| pRG-tk untreated (final) | 12,816 | |
| pRG-tk TPA treated (final) | 11,347 | 0.88 |

Example 3

Synthetic Firefly Luciferase Genes

The *luc+* gene (U.S. Patent No. 5,670,356) was optimized using two approaches. In the first approach (Strategy A), regulatory sequences such as codons were optimized and consensus transcription factor binding sites (TFBS) were removed (see Example 4, although different versions of programs and databases were used). The sequences obtained for the first approach include *hluc+ver2AF1* through *hluc+ver2AF8* (designations with an "F" indicate the construct included flanking sequences). *hluc+ver2AF1* is codon-optimized, *hluc+ver2AF2* is a sequence obtained after a first round of removal of identified undesired sequences including transcription factor binding sites, *hluc+ver2AF3* was obtained after a second round of removal of identified undesired sequences including transcription factor binding sites, *hluc+ver2AF4* was obtained after a third round of removal of identified undesired sequences including transcription factor binding sites, *hluc+ver2AF5* was obtained after a fourth round of removal of identified undesired sequences including transcription factor binding sites, *hluc+ver2AF6* was obtained after removal of promoter modules and RBS, *hluc+ver2AF7* was obtained after further removal of identified undesired sequences including transcription factor binding sites, and *hluc+ver2AF8* was obtained after modifying a restriction enzyme recognition site.

Pairwise DNA identity of different *P.pyralis* luciferase gene versions:

Table 8

| | luc | luc+ | hluc+ | hluc+ver2A1 | hluc+ver2B1 | hluc+ver2A6 | hluc+ver2B6 |
|-------------|-----|------|-------|-------------|-------------|-------------|-------------|
| Luc | 100 | 95 | 76 | 73 | 77 | 74 | 75 |
| luc+ | | 100 | 78 | 76 | 78 | 75 | 77 |
| hluc+ | | | 100 | 91 | 81 | 87 | 81 |
| hluc+ver2A1 | | | | 100 | 74 | 91 | 78 |
| hluc+ver2B1 | | | | | 100 | 74 | 85 |
| hluc+ver2A6 | | | | | | 100 | 80 |
| hluc+ver2B6 | | | | | | | 100 |

luc+ has the following sequence:

atggaagacgcaaaaacataaagaaaggcccgccattctatccgctggaagatggaaccgctggagagca
 actgcataaggctatgaagagatacgccctggctcctggaacaattgctttacagatgcacatatcagggtggacatc
 5 acttacgctgagctacttcgaaatgtccgttcggttggcagaagctatgaaacgatatgggctgaatacaaatcacaga
 atcgtcgtatgcagtgaaaactctctcaattctttatgccggtgttggcgcggtatttatcggagtgcagttgcgccc
 gcgaacgacatttataatgaacgtgaattgtcaacagtatgggcatttcgcagcctaccgtggtgttcgtttccaaa
 aggggttgcaaaaaatttgaacgtgcaaaaaagctcccaatcatcaaaaaattattatcaggattctaaaacgga
 ttaccagggatttcagtcgatgtacacgttcgtcacatctcatctacctcccggttttaataacgattttgtccaga
 10 gtccctcgtatagggacaagacaattgcactgatcatgaactcctctggatctactggtctgcctaaagggtgcgtctg
 cctcatagaactgcctgcgtgagattctcgcagatccagagatcctattttggcaatcaaatcattccggatactgcgat
 ttttaagtgtgttcattccatcacggttttgaatgtttactacactcggatattgatgtggatttcgagtcgtctaat
 gtatagatttgaagaagagctgtttctgaggagccttcaggattacaagattcaaatgcgctgctggtgccaaccta
 ttctccttcttcgcaaaaagcactctgattgacaaatacatttatctaattacacgaaattgcttctggtggcgctcccc
 15 tctctaaggaagtcggggaagcgggttgccaagagggtccatctgccaggtatcaggcaaggatatgggtcactga
 gactacatcagctattctgattacacccgagggggatgataaaccggcgcggtcggttaaagtgttcatttttgaa
 gcgaagggtgtggatctggataccgggaaaacgctggcggttaatacaagaggcgaactgtgtgtgagaggtccta
 tgattatgtccggttatgtaacaatccggaagcgaccaacgccttgattgacaaggatggatggctacattctggag
 acatagcttactgggacgaagacgaacacttctcatcgttgaccgcctgaagtctctgattaagtacaaaggctatca
 20 ggtggctcccgtgaattggaatccatcttgcctcaacaccccaacatcttcgacgcaggtgtcgcaggtcttcccga
 cgatgacgccgggtgaactcccgccgccgttgtgtttggagcacggaaagacgatgacggaaaaagagatcgtg
 gattacgtcgccagtcaagtaacaaccgcgaaaaagtgcgcggaggagttgtgttggacgaagtaccgaaag
 gtcttaccggaaaaactcgacgcaagaaaaatcagagagatcctcataaaggccaagaaggcggaagatcgcc
 gtgtaa (SEQ ID NO:43)

25

and hluc+ has the following sequence:

atggccgatgctaagaacattaagaaggccctgctcccttctaccctctggaggatggcaccgctggcgagcagc
 tgcacaaggccatgaagaggtatgccctggcctggcaccattgccttcaccgatgccacattgaggtggacatc
 acctatgccgagtacttcgagatgtctgtgcgcctggccgaggccatgaagaggtacggcctgaacaccaaccacc
 30 gcatcgtggtgtgctctgagaactctctgcagttcttcatgccagtgtggcgccctgttcacggagtggccgtgg
 cccctgctaacgacatttacaacgagcgcgagctgtgaacagcatgggcatttctcagcctaccgtggtgtcgtgt
 ctaagaaggccctgcagaagatcctgaacgtgcagaagaagctgcctatcatccagaagatcatcatcatggactct
 aagaccgactaccagggttcagagcatgtacacattcgtgacatctcatctgcctcctggcttcaacgagtacgac

ttcgtgccagagtctttcgacagggacaaaaccattgccctgatcatgaacagctctgggtctaccggcctgcctaag
 ggcgtggccctgcctcatcgaccgcctgtgtgcttctctcacgcccgcaccctattttcggcaaccagatcatc
 cccgacaccgctattctgagcgtggtgccattccaccacggcttcggcatgttcaccaccctgggctacctgattgc
 ggctttcgggtggtgctgatgtaccgcttcgaggaggagctgttctgcgcagcctgcaagactacaaaattcagtct
 5 ggcctgtggtgccaaccctgttcagcttctcgctaagagcaccctgatcgacaagtacgacctgttaacctgcac
 gagattgcctctggcggcgccccactgtctaaggaggtgggcgaagccgtggccaagcgctttcatctgccaggca
 tccgccagggtacggcctgaccgagacaaccagcgccattctgattacccagaggcgacgacaagcctggc
 gccgtgggcaaggtggtgccattcttcgaggccaaggtggtggacctggacaccggcaagaccctgggagtga
 ccagcgcggcgagctgtgtgtgcgcccctatgattatgtccggctacgtgaataaccctgaggccacaaacgcc
 10 ctgatcgacaaggacggctggctgcactctggcgacattgcctactgggacgaggacgagcacttctcatctgtgga
 ccgcctgaagtctctgatcaagtacaagggtaccaggtggccccagccgagctggagtctatcctgtgcagcac
 cctaacattttcgacgccggagtgccggcctgcccgaacgacgatgccggcgagctgcctgccgccgtcgtcgtg
 ctggaacacggcaagaccatgaccgagaaggagatcgtggactatgtggccagccaggtgacaaccgccaagaa
 gctgcgcggcgagtggtgttcgtggacgaggtgcccaagggcctgaccggcaagctggacgcccgcaagatcc
 15 gcgagatcctgatcaaggctaagaaaggcggaagatgccgtgtaa (SEQ ID NO:14).

Table 9

Percent Identity

| | | | | | |
|----|-------------------|--------------|--------------|------|-------|
| 20 | | hluc+ver2A8 | hluc+ver2B10 | luc+ | hluc+ |
| | | hluc+ver2A8 | 79.6 | 74 | 86.6 |
| | <u>Divergence</u> | hluc+ver2B10 | 22.9 | 75.9 | 80.1 |
| | | luc+ | 30.4 | 27.8 | 77.4 |
| | | hluc+ | 14.7 | 22.5 | 25.7 |

Table 10

25 Composition statistics of different *P.pyralis* luciferase gene versions

| | GC content | CG di-nucleotides |
|-------------|------------|-------------------|
| H. sapiens | 53% | -- |
| luc | 45% | 99 |
| luc+ | 47% | 97 |
| hluc+ | 60% | 111 |
| hluc+ver2A1 | 66% | 151 |
| hluc+ver2B1 | 46% | 1 |
| hluc+ver2A6 | 58% | 133 |
| hluc+ver2B6 | 49% | 53 |

hluc+ver2A1-hluc+ver2A5 have the following sequences (SEQ ID Nos.16-20):

hluc+ver2A1

5 AAAGCCACCATGGAGGACGCCAAGAACATCAAGAAGGGCCCCGCCC
CCTTCTACCCCCTGGAGGACGGCACCGCCGGCGAGCAGCTGCACAAG
GCCATGAAGCGCTACGCCCTGGTGGCCGGCACCATCGCCTTCACCGA
CGCCACATCGAGGTGGACATCACCTACGCCGAGTACTTCGAGATGA
GCGTGCGCCTGGCCGAGGCCATGAAGCGCTACGGCCTGAACACCAAC
10 CACCGCATCGTGGTGTGCAGCGAGAACAGCCTGCAGTTCTTCATGCC
CGTGCTGGGCGCCCTGTTTCATCGGCGTGGCCGTGGCCCCCGCCAACG
ACATCTACAACGAGCGCGAGCTGCTGAACAGCATGGGCATCAGCCAG
CCCACCGTGGTGTTCGTGAGCAAGAAGGGCCTGCAGAAGATCCTGAA
CGTGCAAGAAGAAGCTGCCCATCATCCAGAAGATCATCATCATGGACA
15 GCAAGACCGACTACCAGGGCTTCCAGAGCATGTACACCTTCGTGACC
AGCCACCTGCCCCCGGCTTCAACGAGTACGACTTCGTGCCCCGAGAG
CTTCGACCGCGACAAGACCATCGCCCTGATCATGAACAGCAGCGGCA
GCACCGGCCTGCCCAAGGGCGTGGCCCTGCCCCACCGCACCGCCTGC
GTGCGCTTCAGCCACGCCCCGCGACCCCATCTTCGGCAACCAGATCAT
20 CCCCACACCGCCATCCTGAGCGTGGTGGCCTTCCACCACGGCTTCG
GCATGTTCAACCACCCTGGGCTACCTGATCTGCGGCTTCCGCGTGGTGC
TGATGTACCGCTTCGAGGAGGAGCTGTTCTGCGCAGCCTGCAGGAC
TACAAGATCCAGAGCGCCCTGCTGGTGGCCACCCTGTTTCAGCTTCTTC
GCCAAGAGCACCTGATCGACAAGTACGACCTGAGCAACCTGCACGA
25 GATCGCCAGCGGCGGCGCCCCCTGAGCAAGGAGGTGGGCGAGGCC
GTGGCCAAGCGCTTCCACCTGCCCCGGCATCCGCCAGGGCTACGGCCT
GACCGAGACCACCAGCGCCATCCTGATCACCCCCGAGGGCGACGACA
AGCCCGGCGCCGTGGGCAAGGTGGTGGCCTTCTTCGAGGCCAAGGTG
GTGGACCTGGACACCGGCAAGACCCTGGGCGTGAACCAGCGCGGCG
30 AGCTGTGCGTGCGCGGCCCCATGATCATGAGCGGCTACGTGAACAAC
CCCGAGGCCACCAACGCCCTGATCGACAAGGACGGCTGGCTGCACAG
CGGCGACATCGCCTACTGGGACGAGGACGAGCACTTCTTCATCGTGG
ACCGCCTGAAGAGCCTGATCAAGTACAAGGGCTACCAGGTGGCCCCC
GCCGAGCTGGAGAGCATCCTGCTGCAGCACCCCAACATCTTCGACGC

CGGCGTGGCCGGCCTGCCCCGACGACGACGCCGGCGAGCTGCCCCGCCG
CCGTGGTGGTGCTGGAGCACGGCAAGACCATGACCGAGAAGGAGAT
CGTGGA~~C~~TACGTGGCCAGCCAGGTGACCACCGCCAAGAAGCTGCGCG
GCGGCGTGGTGTTTCGTGGACGAGGTGCCCAAGGGCCTGACCGGCAAG
5 CTGGACGCCCGCAAGATCCGCGAGATCCTGATCAAGGCCAAGAAGG
GCGGCAAGATCGCCGTGTAATAATTCTAGA

hluc+ver2.A2

AAAGCCACCATGGAGGACGCCAAGAACATCAAGAAGGGGCCAGCGC
10 CATTCTACCCCCTGGAGGACGGCACCGCCGGCGAGCAGCTGCACAAG
GCCATGAAGCGCTACGCCCTGGTGCCCGGCACCATCGCCTTCACCGA
CGCACATATCGAGGTGGACATCACCTACGCCGAGTACTTCGAGATGA
GCGTTCGGCTGGCAGAGGCTATGAAGCGCTATGGGCTGAACACCAAC
CATCGCATCGTGGTGTGCAGCGAGAACAGCTTGCAGTTCTTCATGCC
15 CGTGTTGGGTGCCCTGTTTCATCGGCGTGGCTGTGGCCCCAGCTAACG
ACATCTACAACGAGCGCGAGCTGCTGAACAGCATGGGCATCAGCCAG
CCCACCGTCGTATTCGTGAGCAAGAAAGGGCTGCAAAAGATCCTGAA
CGTGCAAAAGAAGCTGCCCATCATCCAAAAGATCATCATCATGGACA
GCAAGACCGACTACCAGGGCTTCCAAAGCATGTACACCTTCGTGACC
20 AGCCATTTGCCGCCCCGGCTTCAACGAGTACGACTTCGTGCCCGAGAG
CTTCGACCGCGACAAGACCATCGCCCTGATCATGAACAGTAGTGGA
GTACCGGCTTACCTAAGGGCGTGGCCCTACCGCACCGCACCGCCTGT
GTCCGATTTCAGTCATGCCCCGCGACCCCATCTTCGGCAACCAGATCATC
CCCGACACCGCTATCCTGAGCGTGGTGCCATTTACCACGGCTTCGGC
25 ATGTTACCAACCCCTGGGCTACTTGATCTGCGGCTTCCGGGTCGTGCTG
ATGTACCGCTTCGAGGAGGAGCTATTCTTGCGCAGCTTGCAAGACTA
CAAGATTCAAAGCGCCCTGCTGGTGCCACCCCTGTTAGTTTCTTCGC
CAAGAGCACCTGATCGACAAGTACGACCTGAGCAACCTGCACGAG
ATCGCCAGCGGCGGCGCCCCGCTCAGCAAGGAGGTGGGCGAGGCCG
30 TGGCCAAGCGCTTCCACCTGCCAGGCATCCGCCAGGGCTACGGCCTG
ACCGAGACAACCAGCGCCATTCTGATCACCCCCGAGGGGGACGACA
AGCCTGCGCAGTAGGCAAGGTGGTGCCCTTCTTCGAGGCTAAGGTG
GTGGACCTGGACACCGGTAAAACCCTGGGTGTGAACCAGCGCGGCG

AGCTGTGCGTCCGTGGCCCCATGATCATGAGCGGCTACGTTAACAAC
CCCGAGGCTACAAAACGCCCTGATCGACAAGGACGGCTGGCTGCACAG
CGGCGACATCGCCTACTGGGACGAGGACGAGCACTTCTTCATCGTGG
ACCGGCTGAAGAGCCTGATCAAATACAAGGGCTACCAGGTAGCCCCA
5 GCCGAAGTGGAGAGCATCCTGCTGCAGCACCCCAACATCTTCGACGC
CGGGGTCGCCGGCCTGCCCCGACGACGATGCCGGCGAGCTGCCCGCCG
CAGTCGTGGTGCTGGAGCACGGTAAAACCATGACCGAGAAGGAGAT
CGTGGAATATGTGGCCAGCCAGGTTACAACCGCCAAGAAGCTGCGCG
GCGGCGTGGTGTTCTGTGGACGAGGTGCCTAAAGGCCTGACGGGCAAG
10 TTGGACGCCCGCAAAGATCCGCGAGATTCTGATCAAGGCCAAGAAGGG
CGGCAAGATCGCCGTGTAATAATTCTAGA

hluc+ver2A3

AAAGCCACCATGGAAGATGCCAAAAACATTAAGAAGGGGCCAGCGC
15 CATTCTACCCACTGGAGGACGGCACCGCCGGCGAGCAGCTGCACAAA
GCCATGAAGCGCTACGCCCTGGTGCCCGGCACCATCGCCTTTACCGA
CGCACATATCGAGGTGGACATCACCTACGCCGAGTACTTCGAGATGA
GCGTTCGGCTGGCAGAGGCTATGAAGCGCTATGGGCTGAATACCAAC
CATCGCATCGTGGTGTGCAGCGAGAATAGCTTGCAGTTCTTCATGCCC
20 GTGTTGGGTGCCCTGTTTCATCGGTGTGGCTGTGGCCCCAGCTAACGAC
ATCTACAACGAGCGCGAGCTGCTGAACAGCATGGGCATCAGCCAGCC
CACCGTCGTATTCGTGAGCAAGAAAGGGCTGCAAAAGATCCTCAACG
TGCAAAAGAAGCTACCGATCATACAAAAGATCATCATCATGGATAGC
AAGACCGACTACCAAGGGCTTCCAAAGCATGTACACCTTCGTGACCAG
25 CCATTTGCCACCCGCTTCAACGAGTACGACTTCGTGCCCCGAGAGCTT
CGACCGGGACAAAACCATCGCCCTGATCATGAACAGTAGTGGCAGTA
CCGGATTGCCCAAGGGCGTAGCCCTACCGCACCGCACCGCCTGTGTC
CGATTCAATCATGCCCCGCGACCCCATCTTCGGCAACCAGATCATCCCC
GACACCGCTATCCTCAGCGTGGTGCCATTTACACCGGCTTCGGCATG
30 TTCACCACGCTGGGCTACTTGATCTGCGGCTTTCGGGTCGTGCTCATG
TACCGCTTCGAGGAGGAGCTATTCTTGCGCAGCTTGCAAGACTATAA
GATTCAAAGCGCCCTGCTGGTGCCCACTGTTCAGCTTCTTCGCCAA
GAGCACTCTCATCGACAAGTACGACCTGAGCAACCTGCACGAGATCG

CCAGCGGCGGGGCGCCGCTCAGCAA.GGAGGTGGGCGAGGCCGTGGC
CAAGCGCTTCCACCTACCAGGCATCCGCCAGGGCTACGGCCTGACAG
AAACAACCAGCGCCATTCTGATCACCCCCGAAGGGGACGACAAGCCT
GGCGCAGTAGGCAAGGTGGTGCCCTTCTTCGAGGCTAAGGTGGTGGGA
5 CTTGGACACCGGTAAGACCCTGGGTGTGAACCAGCGCGGCGAGCTGT
GCGTCCGTGGCCCCATGATCATGAGCGGCTACGTTAACAACCCCGAG
GCTACAAACGCTCTCATCGACAAGGACGGCTGGCTGCACAGCGGCGA
CATCGCCTACTGGGACGAGGACGAGCACTTCTTCATCGTGGACCGGC
TGAAGAGCCTGATCAAATAACAAGGGCTACCAGGTAGCCCCAGCCGA
10 ACTGGAGAGCATCCTGCTGCAACACCCCAACATCTTCGACGCCGGGG
TCGCCGGCCTGCCCCGACGACGATGCCGGCGAGCTGCCCCGCCGCAGTC
GTCGTGCTGGAGCACGGTAAAACCA.TGACCGAGAAGGAGATCGTGG
ACTATGTGGCCAGCCAGGTTACAACCGCCAAGAAGCTGCGCGGTGGT
GTTGTGTTTCGTGGACGAGGTGCCTAAAGGCCTGACGGGCAAGTTGGA
15 CGCCCGCAAGATCCGCGAGATTCTC.ATTAAGGCCAAGAAGGGCGGCA
AGATCGCCGTGTAATAATTCTAGA

hluc+ver2A4

AAAGCCACCATGGAAGATGCCAAAAACATTAAGAAGGGCCCAGCGC
20 CATTCTACCCACTCGAAGACGGCACCGCCGGCGAGCAGCTGCACAAA
GCCATGAAGCGCTACGCCCTGGTGC.CCGGCACCATCGCCTTTACCGA
CGCACATATCGAGGTGGACATTACCTACGCCGAGTACTTCGAGATGA
GCGTTCGGCTGGCAGAAGCTATGAA.GCGCTATGGGCTGAACACCAAC
CATCGCATCGTGGTGTGCAGCGAGA.ATAGCTTGCA GTTCTTCATGCCC
25 GTGTTGGGTGCCCTGTTTCATCGGTGTGGCTGTGGCCCCAGCTAACGAC
ATCTACAACGAGCGCGAGCTGCTGAACAGCATGGGCATCAGCCAGCC
CACCGTCGTATTTCGTGAGCAAGAAA.GGGCTGCAAAAAGATCCTCAACG
TGCAAAAAGAAGCTACCGATCATACAAAAGATCATCATCATGGATAGC
AAGACCGACTACCAGGGCTTCCAAAGCATGTACACCTTCGTGACTTC
30 CCATTTGCCACCCGGCTTCAACGAGTACGACTTCGTGCCCCGAGAGCTT
CGACCGGGACAAAACCATCGCCCTGATCATGAACAGTAGTGGCAGTA
CCGGATTGCCCAAGGGCGTAGCCCTACCGCACCGCACCGCTTGTGTC
CGATTCAGTCATGCCCGCGACCCCA.TCTTCGGCAACCAGATCATCCCC

GACACCGCTATCCTCAGCGTGGTGCCATTTACCAACGGCTTCGGCATG
TTCACCACGCTGGGCTACTTGATCTGCGGCTTTCGGGTCGTGCTCATG
TACCGCTTCGAGGAGGAGCTATTCTTGCGCAGCTTGCAAGACTATAA
GATTCAAAGCGCCCTGCTGGTGCCACACTGTTCAAGTTTCTTCGCCAA
5 GAGCACTCTCATCGACAAGTACGACCTAAGCAACTTGCACGAGATCG
CCAGCGGCGGGGCGCCGCTCAGCAAGGAGGTGGGCGAGGCCGTGGC
CAAACGCTTCCACCTACCAGGCATCCGCCAGGGCTACGGCCTGACAG
AAACAACCAGCGCCATTCTGATCACCCCCGAAGGGGACGACAAGCCT
GGCGCAGTAGGCAAGGTGGTGCCCTTCTTCGAGGCTAAGGTGGTGGA
10 CTTGGACACCGGTAAGACACTGGGTGTGAACCAGCGCGGCGAGCTGT
GCGTCCGTGGCCCCATGATCATGAGCGGCTACGTTAACAACCCCGAG
GCTACAAACGCTCTCATCGACAAGGACGGCTGGCTGCACAGCGGCGA
CATCGCCTACTGGGACGAGGACGAGCACTTCTTCAATCGTGGACCGGC
TGAAGAGCCTGATCAAATACAAGGGCTACCAGGTAGCCCCAGCCGA
15 ACTGGAGAGCATCCTGCTGCAACACCCCAACATCTTCGACGCCGGGG
TCGCCGGCCTGCCCGACGACGATGCCGGCGAGCTGCCCGCCGAGTC
GTCGTGCTGGAACACGGTAAAACCATGACCGAGAAGGAGATCGTGG
ACTATGTGGCCAGCCAGGTTACAACCGCCAAGAAAGCTGCGCGGTGGT
GTTGTGTTTCGTGGACGAGGTGCCTAAAGGCCTGACGGGCAAGTTGGA
20 CGCCCGCAAGATCCGCGAGATTCTCATTAAAGGCCAAGAAGGGCGGCA
AGATCGCCGTGTAATAATTCTAGA

hluc+ver2A5

AAAGCCACCATGGAAGATGCCAAAAACATTAAGAAAGGGCCCAGCGC
25 CATTCTACCCACTCGAAGACGGCACCGCCGGCGAAGCAGCTGCACAAA
GCCATGAAGCGCTACGCCCTGGTGCCCGGCACCAATCGCCTTTACCGA
CGCACATATCGAGGTGGACATTACCTACGCCGAGTACTTCGAGATGA
GCGTTCGGCTGGCAGAAGCTATGAAGCGCTATGGGCTGAACACCAAC
CATCGGATCGTGGTGTGCAGCGAGAATAGCTTGCAAGTTCTTCATGCC
30 CGTGTTGGGTGCCCTGTTTCATCGGTGTGGCTGTGGCCCCAGCTAACGA
CATCTACAACGAGCGCGAGCTGCTGAACAGCATGGGCATCAGCCAGC
CCACCGTCGTATTCGTGAGCAAGAAAGGGCTGCAAAAGATCCTCAAC
GTGCAAAAGAAGCTACCGATCATAAAAAGATCATCATCATGGATAG

CAAGACCGACTACCAGGGCTTCCAAAGCATGTACACCTTCGTGACTT
 CCCATTTGCCACCCGGCTTCAACGAGTACGACTTCGTGCCCCGAGAGC
 TTCGACCGGGACAAAACCATCGCCCTGATCATGAACAGTAGTGGCAAG
 TACCGGATTGCCCAAGGGCGTAGCCCTACCGCACCGCACCGCTTGTG
 5 TCCGATTTCAGTCATGCCCCGCGACCCCATCTTCGGCAACCAGATCATCC
 CCGACACCGCTATCCTCAGCGTGGTGCCATTTACCACGGCTTCGGCA
 TGTTCACCACGCTGGGCTACTTGATCTGCGGCTTTCGGGTCGTGCTCA
 TGTACCGCTTCGAGGAGGAGCTATTCTTGCGCAGCTTGCAAGACTAT
 AAGATTCAAAGCGCCCTGCTGGTGCCCACTGTTTCAGTTTCTTCGCT
 10 AAGAGCACTCTCATCGACAAGTACGACCTAAGCAACTTGCACGAGAT
 CGCCAGCGGCGGGGCGCCGCTCAGCAAGGAGGTGGGCGAGGCCGTG
 GCCAAACGCTTCCACCTACCAGGCATCCGCCAGGGCTACGGCCTGAC
 AGAAACAACCAGCGCCATTCTGATCACCCCGAAGGGGACGACAAG
 CCTGGCGCAGTAGGCAAGGTGGTGCCCTTCTTCGAGGCTAAGGTGGT
 15 GGACTTGGACACCGGTAAGACACTGGGTGTGAACCAGCGCGGCGAG
 CTGTGCGTCCGTGGCCCCATGATCATGAGCGGCTACGTTAACAACCC
 CGAGGCTACAAACGCTCTCATCGACAAGGACGGCTGGCTGCACAGCG
 GCGACATCGCCTACTGGGACGAGGACGAGCACTTCTTCATCGTGGAAC
 CGGCTGAAGAGCCTGATCAAATACAAGGGCTACCAGGTAGCCCCAGC
 20 CGAACTGGAGAGCATCCTGCTGCAACACCCCAACATCTTCGACGCCG
 GGGTCGCCGGCCTGCCCGACGACGATGCCGGCGAGCTGCCCGCCGCA
 GTCGTCGTGCTGGAACACGGTAAAACCATGACCGAGAAGGAGATCGT
 GGACTATGTGGCCAGCCAGGTTACAACCGCCAAGAAGCTGCGCGGTG
 GTGTTGTGTTTCGTGGACGAGGTGCCTAAAGGCCTGACGGGCAAGTTG
 25 GACGCCCCGAAGATCCGCGAGATTCTCATTAAAGGCCAAGAAGGGCG
 GCAAGATCGCCGTGTAATAATTCTAGA

hluc+ver2A6 has the following sequence

AAAGCCACCATGGAAaGAtGCCAAaAACATtAAGAAGGGCCCaGCgCCaT
 30 TCTACCCaCTcGAaGACGGCACCGCCGGCGAGCAGCTGCACAAaGCCA
 TGAAGCGCTACGCCCTGGTGCCCGGCACCATCGCCTTtACCGACGCaC
 AtATCGAGGTGGACATtACCTACGCCGAGTACTTCGAGATGAGCGTtCG
 gCTGGCaGAaGCtATGAAGCGCTAtGGgCTGAAtACaAACCAAtCGgATCGT

GGTGTGCAGCGAGAAAtAGCtTGCAGTTCTTCATGCCCCGTGtGGGtGCC
 CTGTTTCATCGGtGTGGCtGTGGCCCCaGctAACGACATCTACAACGAGC
 GCGAGCTGCTGAACAGCATGGGCATCAGCCAGCCCACCGTcGTaTTCG
 TGAGCAAGAAaGGgCTGCAaAAGATCCTcAACGTGCAaAAGAAGCTaCC
 5 gATCATaCAaAAGATCATCATCATGGAtAGCAAGACCGACTACCAGGG
 CTTCCAaAGCATGTACACCTTTCGTGACtCcCAAtTGCCaCCCGGCTTCAA
 CGAGTACGACTTCGTGCCCCGAGAGCTTCGACCGgGACAAaACCATCGC
 CCTGATCATGAACAGtAGtGGCAGtACCGGatTgCCcAAGGGCGTgGCCC
 TaCCgCACCGCACCGCtTGTGTcCGaTTCAGtCAtGCCCCGCGACCCCATCT
 10 TCGGCAACCAGATCATCCCCGACACCGCtATCCTcAGCGTGGTGCCaTT
 tCACACGGCTTCGGCATGTTACCCACgCTGGGCTACtTGATCTGCGGC
 TTtCGgGTcGTGCTcATGTACCGCTTCGAGGAGGAGCTaTTCtTGCGCAG
 CtTGCAaGACTAtAAGATtCAaAGCGCCCTGCTGGTGCCCaCTGTTCA
 GtTTCTTCGtAAGAGCACtCTcATCGACAAGTACGACCTaAGCAACtTG
 15 CACGAGATCGCCAGCGGCGGgGCgCCgCTcAGCAAGGAGGTaGGtGAG
 GCCGTGGCCAAaCGCTTCCACCTaCCaGGCATCCGCCAGGGCTACGGC
 CTGACaGAaACaACCAGCGCCATtCTGATCACCCCCGAaGGgGACGACA
 AGCCtGGCGCaGTaGGCAAGGTGGTGCCCTTCTTCGAGGCtAAGGTGGT
 GGACtTGGACACCGGtAAgACaCTGGGtGTGAACCAGCGCGGCGAGCTG
 20 TGCGTcCGtGGCCCCATGATCATGAGCGGCTACGTtAACAACCCCGAG
 GCtACaAACGCtCTcATCGACAAGGACGGCTGGCTGCACAGCGGCGAC
 ATCGCCTACTGGGACGAGGACGAGCACTTCTTCATCGTGGACCGgCT
 GAAGAGCCTGATCAAaTACAAGGGCTACCAGGTaGCCCCaGCCGAaCT
 GGAGAGCATCCTGCTGCAaCACCCCAACATCTTCGACGCCGGgGTcGC
 25 CGGCCTGCCCCGACGACGAtGCCGGCGAGCTGCCCCGCCGCaGTcGTcGT
 GCTGGAaCACGGtAAaACCATGACCGAGAAGGAGATCGTGGACTAtGT
 GGCCAGCCAGGTtACaACCGCCAAGAAGCTGCGCGGtGGtGTtGTGTTC
 GTGGACGAGGTGCCtAAaGGCCTGACgGGCAAGtTGGACGCCCGCAAG
 ATCCGCGAGATtCTcATtAAGGCCAAGAAGGGCGGCAAGATCGCCGTG
 30 TAATAATTCTAGA (SEQ ID NO:21).

The hluc+ver2A6 sequence was modified yielding hluc+ver2A7:

AAAGCCACCATGGAaGAtGCCAAaAACATtAAGAA
GGGCCCgGCgCCaTTCTACCCaCTcGAAgACGGgACCGCCGGCGAGCAG
CTGCACAAaGCCATGAAGCGCTACGCCCTGGTGGCCGGCACCATCGC
CTTtACCGACGCaCAtATCGAGGTGGACATtACCTACGCCGAGTACTTC
5 GAGATGAGCGTtCGgCTGGCaGAaGCTATGAAGCGCTAtGGgCTGAAtAC
aAACCAAtCGgATCGTGGTGTGCAGCGAGAAAtAGCtTGCAGTTCTTCATGC
CCGTGtTGGGtGCCCTGTTcATCGGtGTGGCtGTGGCCCCaGCtAACGAC
ATCTACAACGAGCGCGAGCTGCTGAACAGCATGGGCATCAGCCAGCC
CACCGTcGTaTTCGTGAGCAAGAAaGGgCTGCAaAAGATCCTcAACGTG
10 CAaAAGAAGCTaCCgATCATaCAaAAGATCATCATCATGGAtAGCAAGA
CCGACTACCAGGGCTTCCaAAGCATGTACACCTTCGTGACtCcCAAtTG
CCaCCCGGCTTCAACGAGTACGACTTCGTGCCCGAGAGCTTCGACCGg
GACAAaACCATCGCCCTGATCATGAACAGtAGtGGCAGtACCGGAtTgCC
cAAGGGCGTaGCCCTaCCgCACCGCACCGCtTGtGTcCGaTTCAGtCAtGCC
15 CGCGACCCCATCTTCGGCAACCAGATCATCCCCGACACCGCtATCCTc
AGCGTGGTGCCaTTtCACACGGCTTCGGCATGTTACCCACgCTGGGCT
ACtTGATCTGCGGCTTtCGgGTcGTGCTcATGTACCGCTTCGAGGAGGAG
CTaTTCtTGCGCAGCtTGCAaGACTAtAAGATtCAatctGCCCTGCTGGTGC
CCACaCTaTtTAGcTTCTTCGCTAAGAGCACTcTcATCGACAAGTACGACC
20 TaAGCAACtTGCACGAGATCGCCAGCGGGCGgGCgCCgCTcAGCAAGGA
GGTaGGtGAGGCCGTGGCCAAaCGCTTCCACCTaCCaGGCATCCGCCAG
GGCTACGGCCTGACaGAaACaACCAGCGCCATtCTGATCACCCCGAaG
GgGACGACAAGCCtGGCGCaGTaGGCAAGGTGGTGGCCTTCTTCGAGG
CtAAGGTGGTGGACtTGGACACCGGtAAgACaCTGGGtGTGAACCAGCG
25 CGGCGAGCTGTGCGTcCGtGGCCCCATGATCATGAGCGGCTACGTtAA
CAACCCCGAGGCTACaAACGCtCTcATCGACAAGGACGGCTGGCTGCA
CAGCGGCGACATCGCCTACTGGGACGAGGACGAGCACTTCTTCATCG
TGGACCGgCTGAAGAGCCTGATCAAAaTACAAGGGCTACCAGGTaGCCC
CaGCCGAaCTGGAGAGCATCCTGCTGCAaCACCCCAACATCTTCGACG
30 CCGGgGTcGCCGGCCTGCCCCGACGACGAAtGCCGGCGAGCTGCCCCGG
CaGTcGTcGTGCTGGAaCACGGtAAaACCATGACCGAGAAGGAGATCGT
GGACTAtGTGGCCAGCCAGGTtACaACCGCCAAGAAGCTGCGCGGtGGt
GTtGTGTTCGTGGACGAGGTGCCtAAaGGCCTGACgGGCAAGtTGGACG

CCCGCAAGATCCGCGAGATtCTcATtAAGGCCAAGAAGGGCGGCAAGA
TCGCCGTGTAATAATTCTAGA (SEQ ID NO:22).

For vectors with a *Bgl*I site in the multiple cloning region, the *Bgl*I site present in
5 the firefly sequence can be removed. The luciferase gene from hluc+ver2AF8,
which lacks a *Bgl*I site, displays an average of a 7.2-fold increase in expression
when assayed in four mammalian cell lines, i.e., NIH3T3, CHO, HeLa and
HEK293 cells.

10 hluc+ver2A8 has the following sequence:

AAAGCCACCATGGAaGAtGCCAAaAACATtAAGAAGGGCCCaGcCCaT
TCTACCCaCTcGAaGACGGgACCGCCGGCGAGCAGCTGCACAAaGCCA
TGAAGCGCTACGCCCTGGTGCCCGGCACCATCGCCTTtACCGACGCaC
AtATCGAGGTGGACATtACCTACGCCGAGTACTTCGAGATGAGCGTtCG
15 gCTGGCaGAaGcAtGAAGCGCTAtGGgCTGAAtACaAACCAtCGgATCGT
GGTGTGCAGCGAGAAAtAGCtTGCAGTTCTTCATGCCCCGTGtTGGGtGCC
CTGTTCATCGGtGTGGCtGTGGCCCCaGcTAACGACATCTACAACGAGC
GCGAGCTGCTGAACAGCATGGGCATCAGCCAGCCCACCGTcGTaTTCCG
TGAGCAAGAAaGGgCTGCAaAAGATCCTcAACGTGCAaAAGAAGCTaCC
20 gATCATaCAaAAGATCATCATCATGGAtAGCAAGACCGACTACCAGGG
CTTCCAaAGCATGTACACCTTCGTGACttcCCAAtTGCCaCCCGGCTTCAA
CGAGTACGACTTCGTGCCCCGAGAGCTTCGACCGgGACAAaACCATCGC
CCTGATCATGAACAGtAGtGGCAGtACCGGAtTgCCcAAGGGCGTaGCCC
TaCCgCACCGCACCGCtTGtGTcCGaTTCAGtCAtGCCCCGCGACCCCATCT
25 TCGGCAACCAGATCATCCCCGACACCGCtATCCTcAGCGTGGTGCCaTT
tACCACGGCTTCGGCATGTTCACCACgCTGGGCTACtTGATCTGCGGC
TTtCGgGTcGTGCTcATGTACCGCTTCGAGGAGGAGCTaTTCtTGCGCAG
CtTGCAaGACTAtAAGATtCAatctGCCCTGCTGGTGCCCaCTaTTtAGcT
TCTTCGtAAGAGCACtCTcATCGACAAGTACGACCTaAGCAACtTGCAC
30 GAGATCGCCAGCGGGCGGgGCgCCgCTcAGCAAGGAGGTaGGtGAGGCC
GTGGCCAAaCGCTTCCACCTaCCaGGCATCCGCCAGGGCTACGGCCTG
ACaGAaACaACCAGCGCCATtCTGATCACCCCCGAaGGgGACGACAAGC
CtGGCGCaGTaGGCAAGGTGGTGCCCTTCTTCGAGGctAAGGTGGTGGa
CtTGGACACCGGtAAgACaCTGGGtGTGAACCAGCGCGGCGAGCTGTGC

GTcCGtGGCCCCATGATCATGAGCGGCTACGTtAACAACCCCGAGGcTA
 CaAACGCtCTcATCGACAAGGACGGCTGGCTGCACAGCGGCGACATCG
 CCTACTGGGACGAGGACGAGCACTTCTTCATCGTGGACCGgCTGAAG
 AGCCTGATCAAAaTACAAGGGCTACCAGGTaGCCCCaGCCGAaCTGGAG
 5 AGCATCCTGCTGCAaCACCCCAACATCTTCGACGCCGGgGTcGCCGGC
 CTGCCCCGACGACGAtGCCGGCGAGCTGCCCCGCCGCaGTcGTcGTGCTGG
 AaCACGGtAAaACCATGACCGAGAAGGAGATCGTGGACTAtGTGGCCA
 GCCAGGTtACaACCGCCAAGAAGCTGCGCGGtGGtGTtGTGTTCGTGGA
 CGAGGTGCCtAAaGGaCTGACcGGCAAGtTGGACGCCCGCAAGATCCGC
 10 GAGATtCTcATtAAGGCCAAGAAGGGCGGCAAGATCGCCGTGTAATAA
 TTCTAGA (SEQ ID NO:23).

For the second approach, firefly luciferase *luc+* codons were optimized for
 mammalian expression, and the number of consensus transcription factor binding
 15 site, and CG dinucleotides (CG islands, potential methylation sites) was reduced.
 The second approach yielded: versions hluc+ver2BF1 through hluc+ver2BF5.
 hluc+ver2BF1 is codon-optimized, hluc+ver2BF2 is a sequence obtained after a
 first round of removal of identified undesired sequences including transcription
 factor binding sites, hluc+ver2BF3 was obtained after a second round of removal
 20 of identified undesired sequences including transcription factor binding sites,
 hluc+ver2BF4 was obtained after a third round of removal of identified
 undesired sequences including transcription factor binding sites, hluc+ver2BF5
 was obtained after a fourth round of removal of identified undesired sequences
 including transcription factor binding sites, hluc+ver2BF6 was obtained after
 25 removal of promoter modules and RBS, hluc+ver2BF7 was obtained after further
 removal of identified undesired sequences including transcription factor binding
 sites, and hluc+ver2BF8 was obtained after modifying a restriction enzyme
 recognition site.

30 hluc+ver2B1-B5 have the following sequences (SEQ ID Nos. 24-28):
 hluc+ver2B1
 AAAGCCACCATGGAGGATGCTAAGAATATTAAGAAGGGGCCTGCTCC
 TTTTATCCTCTGGAGGATGGGACAGCTGGGGAGCAGCTGCATAAGG

CTATGAAGAGATATGCTCTGGTGCCTGGGACAATTGCTTTTACAGATG
CTCATATTGAGGTGGATATTACATATGCTGAGTATTTTGAGATGTCTG
TGAGACTGGCTGAGGCTATGAAGAGATATGGGCTGAATACAAATCAT
AGAATTGTGGTGTGTTCTGAGAATTCTCTGCAGTTTTTTATGCCTGTG
5 CTGGGGGCTCTGTTTATTGGGGTGGCTGTGGCTCCTGCTAATGATATT
TATAATGAGAGAGAGCTGCTGAATTCTATGGGGATTTCTCAGCCTAC
AGTGGTGTGTTGTGTCTAAGAAGGGGCTGCAGAAGATTCTGAATGTGC
AGAAGAAGCTGCCTATTATTCAGAAGATTATTATTATGGATTCTAAG
ACAGATTATCAGGGGTTTCAGTCTATGTATACATTTGTGACATCTCAT
10 CTGCCTCCTGGGTTTAATGAGTATGATTTTGTGCCTGAGTCTTTTGAT
AGAGATAAGACAATTGCTCTGATTATGAATTCTTCTGGGTCTACAGG
GCTGCCTAAGGGGGTGGCTCTGCCTCATAGAACAGCTTGTGTGAGAT
TTTCTCATGCTAGAGATCCTATTTTTGGGAATCAGATTATTCCTGATA
CAGCTATTCTGTCTGTGGTGCCTTTTCATCATGGGTTTGGGATGTTTAC
15 AACACTGGGGTATCTGATTTGTGGGTTTAGAGTGGTGCTGATGTATAG
ATTTGAGGAGGAGCTGTTTCTGAGATCTCTGCAGGATTATAAGATTCA
GTCTGCTCTGCTGGTGCCTACACTGTTTTCTTTTTTTGCTAAGTCTACA
CTGATTGATAAGTATGATCTGTCTAATCTGCATGAGATTGCTTCTGGG
GGGGCTCCTCTGTCTAAGGAGGTGGGGGAGGCTGTGGCTAAGAGATT
20 TCATCTGCCTGGGATTAGACAGGGGTATGGGCTGACAGAGACAACAT
CTGCTATTCTGATTACACCTGAGGGGGATGATAAGCCTGGGGCTGTG
GGGAAGGTGGTGCCTTTTTTTGAGGCTAAGGTGGTGGATCTGGATAC
AGGGAAGACACTGGGGGTGAATCAGAGAGGGGAGCTGTGTGTGAGA
GGGCCTATGATTATGTCTGGGTATGTGAATAATCCTGAGGCTACAAA
25 TGCTCTGATTGATAAGGATGGGTGGCTGCATTCTGGGGATATTGCTTA
TTGGGATGAGGATGAGCATTTTTTTTATTGTGGATAGACTGAAGTCTCT
GATTAAGTATAAGGGGTATCAGGTGGCTCCTGCTGAGCTGGAGTCTA
TTCTGCTGCAGCATCCTAATATTTTTGATGCTGGGGTGGCTGGGCTGC
CTGATGATGATGCTGGGGAGCTGCCTGCTGCTGTGGTGGTGTGGAG
30 CATGGGAAGACAATGACAGAGAAGGAGATTGTGGATTATGTGGCTTC
TCAGGTGACAACAGCTAAGAAGCTGAGAGGGGGGGTGGTGTGTTGTGG
ATGAGGTGCCTAAGGGGCTGACAGGGAAGCTGGATGCTAGAAAGAT
TAGAGAGATTCTGATTAAGGCTAAGAAGGGGGGGAAGATTGCTGTGT

AATAATTCTAGA

hluc+ver2B2

AAAGCCACCATGGAAGATGCTAAAAACATTAAGAAGGGGCCTGCTCC
5 TTTCTACCCTCTGGAGGATGGGACTGCCGGGAGCAGCTGCATAAAG
CTATGAAGCGGTATGCTCTGGTGCCAGGCACAATTGCGTTCACGGAT
GCTCACATTGAGGTGGACATTACATACGCTGAGTATTTTGAGATGTCG
GTGCGGCTGGCTGAGGCTATGAAGCGATATGGGCTGAATACAAACCA
TAGAATTGTAGTGTGCTCTGAGAACTCGTTGCAGTTTTTTATGCCTGT
10 GCTGGGGGCTCTCTTCATCGGGGTGGCTGTGGCTCCTGCTAACGACAT
TTACAATGAGAGAGAGCTTTTGAACTCGATGGGGATTCTCAGCCTA
CAGTGGTGTTTGTGAGTAAGAAAGGGCTTCAAAGATTCTCAATGTG
CAAAGAAGCTGCCTATTATTCAAAGATTATTATTATGGACTCTAA
GACAGACTACCAGGGGTTTCAGTCTATGTATACATTTGTGACATCTCA
15 TCTGCCTCCTGGGTTCACGAGTATGACTTTGTGCCCCGAGTCTTTCGA
CAGAGATAAGACAATTGCTCTGATTATGAATTCATCTGGGTCTACCG
GGCTGCCTAAGGGTGTAGCTCTGCCACATAGAACAGCTTGTGTGAGA
TTTTCTCATGCTAGGGACCCTATTTTTGGGAATCAGATTATTCCTGAT
ACTGCTATTCTGTTCGGTTGTGCCCTTTCATCATGGGTTTGGGATGTTTA
20 CAACACTGGGCTACCTGATATGTGGGTTTAGAGTGGTGCTCATGTATA
GGTTTGAGGAGGAGCTTTTTTTGCGCTCTCTGCAAGATTATAAGATTC
AGTCTGCTCTGCTGGTGCCTACACTGTTTTCTTTTTTTGCTAAGTCTAC
CCTGATCGATAAGTATGATCTGTCCAACCTGCACGAGATTGCTTCTGG
GGGGGCTCCTCTGTCTAAGGAGGTAGGTGAGGCTGTGGCTAAGCGCT
25 TTCATCTGCCTGGAATCAGACAGGGGTATGGGCTAACAGAAACAACA
TCTGCTATTCTGATTACACCAGAGGGGGATGATAAGCCCGGGGCTGT
AGGGAAAGTGGTGCCCTTTTTTTGAAGCTAAAGTAGTTGATCTTGATAC
CGGTAAGACACTGGGGGTGAATCAGCGAGGGGAAGTGTGTGTGAGA
GGGCCTATGATTATGTGCGGGTATGTGAACAACCCTGAGGCTACAAA
30 TGCTCTGATTGATAAGGATGGGTGGCTGCATTCGGGCGATATTGCTTA
CTGGGATGAGGATGAGCATTTCTTCATCGTGGACAGACTGAAGTCGT
TGATCAAATATAAGGGGTATCAAGTAGCTCCTGCTGAGCTGGAGTCC
ATTCTGCTTCAACATCCTAACATTTTCGATGCTGGGGTGGCTGGGCTG

CCTGATGATGATGCTGGGGAGCTGCCTGCTGCTGTAGTGGTGCTGGA
GCACGGTAAGACAATGACAGAGAAGGAGATTGTGGATTATGTGGCTT
CACAAGTGACAACAGCTAAGAACTGAGAGGTGGCGTTGTGTTTGTG
GATGAGGTGCCTAAAGGGCTGACAGGCAAGCTGGATGCTAGAAAAA
5 TTCGAGAGATTCTGATTAAGGCTAAGAAGGGTGGAAGATTGCTGTG
TAATAGTTC TAGA

hluc+ver2B3

AAAGCCACCATGGAAGATGCTAAAAACATTAAGAAGGGGGCCTGCTCC
10 TTTCTACCCTCTTGAAGATGGGACTGCTGGCGAGCAACTTCACAAAG
CTATGAAGCGGTATGCTCTTGTGCCAGGCACAATTGCGTTCACGGAT
GCTCACATTGAGGTGGACATCACATACGCTGAGTATTTTGAGATGTC
GGTGCGGCTGGCAGAAGCTATGAAGCGCTATGGGCTGAATACAAACC
ATAGAATTGTAGTGTGCAGTGAGAACTCGTTGCAGTTCTTTATGCCCG
15 TGCTGGGGGCTCTCTTCATCGGGGTGGCTGTGGCTCCTGCTAACGACA
TCTACAACGAGCGAGAGCTGTTGAACTCGATGGGGATTTCTCAGCCT
ACAGTGGTGTTTGTGAGTAAGAAAGGGCTTCAAAGATTCTCAATGT
GCAAAAGAAGCTGCCTATTATTCAAAGATTATTATTATGGACTCTA
AGACCGACTACCAGGGGTTTCAGTCTATGTATACATTTGTGACATCTC
20 ATCTGCCTCCTGGCTTCAACGAGTACGACTTCGTGCCCCGAGTCTTTCG
ACAGAGATAAGACAATTGCTCTGATCATGAATTCATCCGGGTCTACC
GGGCTGCCTAAGGGTGTAGCTCTGCCCCATAGAACAGCTTGTGTGAG
ATTTTCTCATGCTAGGGACCCTATTTTGGGAATCAGATTATTCCTGA
CACTGCTATTCTGTGCGGTGGTGCCCTTTCATCATGGGTTTGGGATGTT
25 TACAACACTGGGCTACCTAATATGTGGGTTTAGAGTGGTGCTCATGTA
TAGGTTTGAAGAAGAGCTGTTCTTACGCTCTTTGCAAGATTATAAGAT
TCAGTCTGCTCTGCTGGTGCCAACACTATTCTCTTTTTTGTAAAGTCT
ACGCTCATA GACAAGTATGACTTGTCCAACCTGCACGAGATTGCTTCT
GGCGGAGCACCTCTGTCTAAGGAGGTAGGTGAGGCTGTGGCTAAGCG
30 CTTTCATCTGCCTGGTATCAGACAGGGGTATGGGCTAACAGAAACAA
CATCTGCTATTCTGATTACACCAGAGGGGGATGATAAGCCCCGGGGCT
GTAGGGAAAAGTGGTGCCCTTTTTTGAAGCCAAAGTAGTTGATCTTGAT
ACCGGTAAGACACTAGGGGTGAACCAGCGTGGTGAACCTGTGTGTGAG

AGGGCCTATGATTATGTCGGGGTACGTTAACAACCCCGAAGCTACAA
ATGCTCTGATTGATAAGGATGGCTGGCTGCATTCTGGGCGACATTGCTT
ACTGGGATGAGGATGAGCATTCTTCATCGTGGACAGACTGAAGTCG
TTGATCAAATACAAGGGGTATCAAGTAGCTCCTGCTGAGCTGGAATC
5 CATTCTGCTTCAACATCCC AACATTTTCGATGCTGGGGTGGCTGGGCT
GCCTGATGATGATGCTGGGGAGTTGCCTGCTGCTGTAGTGGTGCTTGA
GCACGGTAAGACAATGACAGAGAAGGAGATCGTGGATTATGTGGCTT
CACAAGTGACAACAGCTAAGAACTGAGAGGTGGCGTTGTGTTTGTG
GATGAGGTGCCTAAAGGGCTCACTGGCAAGCTGGATGCTAGAAAAAT
10 TCGAGAGATTCTGATTAAAGGCTAAGAAGGGTGGAAAGATTGCTGTGT
AATAGTTCTAGA

hluc+ver2B4

AAAGCCACCATGGAAGATGCTAAAAACATTAAGAAGGGGCCTGCTCC
15 CTTCTACCCTCTTGAAGATGGGACTGCTGGCGAGCAACTTCACAAAG
CTATGAAGCGGTATGCTCTTGTGCCAGGCACAATTGCGTTCACGGAT
GCTCACATTGAGGTGGACATCACATACGCTGAGTATTTTGAGATGTC
GGTGC GGCTGGCAGAAGCTATGAAGCGCTATGGGCTGAATACAAACC
ATAGAATTGTAGTGTGCAGTGAGAACTCGTTGCAGTTCTTTATGCCCCG
20 TGCTGGGGGCTCTCTTCATCGGGGTGGCTGTGGCTCCTGCTAACGACA
TCTACAACGAGCGAGAGCTGTTGAACTCGATGGGGATCTCTCAGCCT
ACAGTGGTGTGTTGTGAGTAAGAAAGGGCTTCAAAGATTCTCAATGT
GCAAAGAAGCTGCCTATTATTCAAAGATTATTATTATGGACTCTA
AGACAGACTACCAGGGGTTCAGTCCATGTATACATTTGTGACATCTC
25 ATCTGCCTCCTGGCTTCAACGAGTACGACTTCGTGCCCCGAGTCTTTCG
ACAGAGATAAGACAATTGCTCTGATCATGAATTCATCCGGGTCTACC
GGGCTGCCTAAGGGTGTA GCTCTGCCCCATCGAACAGCTTGTGTGAG
ATTCTCTCATGCCAGGGACCCGATCTTTGGGAATCAGATTATTCCTGA
CACTGCTATTCTGTCGGTGGTGCCCTTTCATCATGGGTTTGGGATGTT
30 TACAACACTGGGATACCTAATATGTGGGTTTAGAGTGGTGCTCATGT
ATAGGTTTGAAGAAGAACTGTTCTTACGCTCTTTGCAAGATTATAAGA
TTCAGTCTGCTCTGCTGGTGCCAACACTATTCTCTTTTTTTGCTAAGTC
TACGCTCATAGACAAGTATGACTTGTCCAACCTGCACGAGATTGCTTC

TGGCGGAGCACCTCTGTCTAAGGAGGTAGGTGAGGCTGTGGCTAAGC
 GCTTTCATCTGCCTGGTATCAGACAGGGGTACGGGCTAACAGAAACA
 ACTTCTGCTATTCTGATTACACCAGAGGGCGATGACAAGCCCCGGGGC
 TGTAGGGAAAGTGGTGGCCCTTTTTTGAAGC CAAAGTAGTTGATCTTGA
 5 TACCGGTAAGACACTAGGGGTGAACCAGCGTGGTGAAGTGTGTGTGC
 GGGGCCCTATGATTATGTCGGGGTACGTTAACAACCCCGAAGCTACA
 AATGCTCTTATTGATAAGGATGGCTGGTTGCATTCGGGCGACATTGCC
 TACTGGGATGAGGATGAGCATTCTTCATC GTGGACAGACTGAAGTC
 GTTGATCAAATACAAGGGGTATCAAGTAGCTCCTGCTGAGCTGGAAT
 10 CCATTCTGCTTCAACATCCAAACATTTTCGATGCTGGGGTGGCTGGGC
 TGCCTGATGATGATGCTGGAGAGTTGCCTGCTGCTGTAGTAGTGCTTG
 AGCACGGTAAGACAATGACAGAGAAGGAGATCGTGGATTATGTGGC
 TTCACAAGTGACAACAGCTAAGAACTGAGAGGTGGCGTTGTGTTTG
 TGGATGAGGTGCCTAAAGGGCTCACTGGCAAGCTGGATGCCAGAAAA
 15 ATTCGAGAGATTCTCATTAAAGGCTAAGAAGGGTGGAAAGATTGCTGT
 GTAATAGTTCTAGA

hluc+ver2B5

AAAGCCACCATGGAAGATGCTAAAAACATTAAGAAGGGGCCTGCTCC
 20 CTTCTACCCTCTTGAAGATGGGACTGCTGGCGAGCAACTTCACAAAG
 CTATGAAGCGGTATGCTCTTGTGCCAGGCACAATTGCGTTCACGGAT
 GCTCACATTGAGGTGGACATCACATACGCTGAGTATTTTGAGATGTC
 GGTGCGGCTGGCAGAAGCTATGAAGCGCTATGGGCTGAATACAAACC
 ATAGAATTGTAGTGTGCAGTGAGAACTCGTTGCAGTTCTTTATGCCCG
 25 TGCTGGGGGCTCTCTTCATCGGGGTGGCTGTGGCTCCTGCTAACGACA
 TCTACAACGAGCGAGAGCTGTTGAACTCGATGGGGATCTCTCAGCCT
 ACAGTGGTGTGTTGTGAGTAAGAAAGGGCTTCAAAAGATTCTCAATGT
 GCAAAAGAAGCTGCCTATTATACAAAAGATTATTATTATGGACTCTA
 AGACCGACTACCAGGGGTTTCAGTCCATGTACACATTTGTAACCTCTC
 30 ATCTGCCTCCTGGCTTCAACGAGTACGACTTCGTGCCCCGAGTCTTTCG
 ACAGGGACAAAACGATTGCTCTGATCATGAACTCATCCGGGTCTACC
 GGGCTGCCTAAGGGTGTAGCTCTGCCCCATCGAACAGCTTGTGTGAG
 ATTCTCTCATGCCAGGGACCCGATCTTTGGGAATCAGATTATTCCTGA

CACTGCTATTCTGTCCGGTGGTGCCCTTTTCATCATGGGTTT GGGATGTT
 CACAACACTGGGATACCTCATTTGCGGGTTTAGAGTGGT GCTCATGTA
 TAGGTTTGAAGAAGAAGTATTCTACGCTCTTTGCAAGATTATAAGAT
 TCAGTCTGCTCTGCTGGTGCCAACTATTCTCTTTTTTTGCTAAGTCT
 5 ACGCTCATAGACAAGTATGACTTGTCCAAGTGCACGAGATTGCTTCT
 GGCGGAGCACCTCTGTCTAAGGAGGTAGGTGAGGCTGTGGCTAAGCG
 CTTTCATCTGCCTGGTATCAGACAGGGGTACGGGCTAACAGAAACAA
 CTTCTGCTATTCTGATTACACCAGAGGGCGATGACAAACCCGGGGCT
 GTAGGGAAAGTGGTGCCCTTTTTTGAAGCCAAAGTAGTTGATCTTGAT
 10 ACCGGTAAGACACTAGGGGTGAACCAGCGTGGTGAAGTGTGTGTGCG
 GGGCCCTATGATTATGTCGGGGTACGTTAACAACCCCGAAGCTACAA
 ATGCTCTTATTGATAAGGATGGCTGGTTGCATTTCGGGCGACATTGCCT
 ACTGGGATGAGGATGAGCATTCTTCATCGTGGACAGACTGAAGTCG
 TTGATCAAATACAAGGGGTATCAAGTAGCTCCTGCTGAGCTGGAATC
 15 CATTCTGCTTCAACATCCTAACATTTTCGATGCTGGGGTGGCTGGGCT
 GCCTGATGATGATGCTGGAGAGTTGCCTGCTGCTGTAGTAGTGCTTGA
 GCACGGTAAGACAATGACAGAGAAGGAGATCGTGGATTATGTGGCTT
 CACAAGTGACAACAGCTAAGAACTGAGAGGTGGCGTTGTGTTTGTG
 GATGAGGTGCCTAAAGGGCTCACTGGCAAGCTGGATGCAGAAAAAT
 20 TCGAGAGATTCTCATTAAAGGCTAAGAAGGGTGGAAAGATTGCTGTGT
 AATAGTTCTAGA

hluc+ver2B6 has the following sequence:

AAAGCCACCATGGAaGATGCcAAaAAcATTAAGAAGGGGCCTGCTCCc
 25 TTcTAcCCTCTtGAaGATGGGACtGcTGGcGAGCAaCTtCAcAAaGCTATGA
 AGcGgTATGCTCTtGTGCCaGGcACAATTGCgTTcACgGATGCTCAcATTG
 AaGTaGAcATcACATAcGCTGAGTATTTTGAGATGTCgGTGcGgCTGGCa
 GAaGCTATGAAGcGcTATGGGCTGAATACAAAcCATAGAATTGTaGTGT
 GcagTGAGAAcTCgtTGCAGTTcTTTATGCCcGTGCTGGGGGCTCTcTTcAT
 30 cGGGGTGGCTGTGGCTCCTGCTAAcGAcATcTAcAAcGAGcGAGAGCTgt
 TGAAcTCgATGGGGATcTCTCAGCCTACAGTGGTGTtGTGagTAAGAA
 aGGGCTtCAaAAGATTCTcAATGTGCAaAAGAAGCTGCCTATTATaCAaA
 AGATTATTATTATGGAcTcTAAGACcGAcTAcCAGGGGTTTCAGTCcATG

TAcACATTTGTaACcTCTCATCTGCCTCCTGGcTTcAAcGAGTAcGAcTTc
 GTGCCcGAGTCTTTcGAcAGgGAcAAaACgATTGCTCTGATcATGAAcagc
 TcGGGTCTACcGGGCTGCCTAAGGGtGTaGCTCTGCCcCATcGAACAGC
 TTGTGTGAGATTcTCTCATGCcAGgGAcCCgATcTTtGGaAAcCAGATcATc
 5 CCTGAcAcTgCTATTCTGTcGtGTgGTGCCcTTTCATCATGGGTTTGGGAT
 GTTcACAACACTGGGATAcTcATtTGcGGGTTTAGAGTGGTGCTcATGTA
 TAGgTTTGAAgAaGAaCTaTTccTacGcTCTtTGCAaGATTATAAGATTcAG
 TCTGCTCTGCTGGTGCCaACACTaTTcTCTTTTTTTGCTAAGTCTACgCTc
 ATaGAcAAGTATGActTGTCCAAcTGCAcGAGATTGCTTCTGGcGGaGCa
 10 CCTCTGTCTAAGGAGGTaGGtGAGGCTGTGGCTAAGcGcTTTCATCTGC
 CTGGtATcAGACAGGGGTAcGGGCTaACAGAAACAACtTCTGCTATTCTG
 ATTACACCAGAGGGcGATGAcAAaCCcGGGGCTGTaGGGAaGTGGTGC
 CcTTTTTTGAaGCcAAaGTaGTtGATCTtGATACcGGtAAGACACTaGGGGT
 GAAcCAGcGtGGtGAaCTGTGTGTGcGgGGcCCTATGATTATGTCgGGGTA
 15 cGTtAAcAAcCCcGAaGCTACAAATGCTCTcATaGAcAAGGAcGGgTGGcTt
 CATagcGGcGAcATTGCcTAcTGGGAcGAGGATGAGCATTTCtTcATcGTG
 GAcAGACTGAAGTCgtTGATcAAaTAcAAGGGGTATCAaGTaGCTCCTGC
 TGAGCTGGAaTcATTCTGCTtCAaCAcCCcAAATcTtGATGCTGGGGT
 GGCTGGGCTGCCTGATGATGATGCTGGaGAGcTGCCTGCTGCTGTaGTa
 20 GTGCTtGAGCAcGGtAAGACAATGACAGAGAAGGAGATcGTGGATTAT
 GTGGCTTCaCAaGTGACAACAGCTAAGAAaCTGAGAGGtGGcGTtGTGT
 TTGTGGATGAGGTGCCTAAaGGGCTcACtGGcAAGCTGGATGCcAGAAA
 aATTcGAGAGATTCTcATTAAAGGCTAAGAAGGGtGGaAAGATTGCTGTG
 TAATAgTTCTAGA (SEQ ID NO:29).

25

hluc+ver2BF8 was created by removing a *Ptx1* consensus transcription factor
 binding site from hluc+ver2BF7.

hluc+ver2B7 has the following sequence:

30 AAAGCCACCATGGAAGATGCCAAAAACATTAAGAAGGGGCCTGCTC
 CCTTCTACCCTCTTGAAGATGGGACTGCTGGCGAGCAACTTCACAAA
 GCTATGAAGCGGTATGCTCTTGTGCCAGGGACAATTGCGTTCACGGA
 TGCTCACATTGAAGTAGACATCACATACGCTGAGTATTTTGAGATGTC

GGTGCGGCTGGCAGAAGCTATGAAGCGCTATGGGCTGAATACAAACC
ATAGAATTGTAGTGTGCAGTGAGAACTCGTTGCAGTTCTTTATGCCCCG
TGCTGGGGGGCTCTCTTCATCGGGGTGGCTGTGGCTCCTGCTAACGACA
TCTACAACGAGCGAGAGCTGTTGAACTCGATGGGGATCTCTCAGCCT
5 ACAGTGGTGTGTTGTGAGTAAGAAAGGGCTTCAAAAGATTCTCAATGT
GCAAAAGAAGCTGCCTATTATACAAAAGATTATTATTATGGACTCTA
AGACAGACTACCAGGGGTTTCAGTCCATGTACACATTTGTAACCTCTC
ATCTGCCTCCTGGCTTCAACGAGTACGACTTCGTGCCCCGAGTCTTTCG
ACAGGGACAAAACGATTGCTCTGATCATGAACAGCTCCGGGTCTACC
10 GGGCTGCCTAAGGGTGTAGCTCTGCCCCATCGAACAGCTTGTGTGAG
ATTCTCTCATGCCAGGGACCCGATCTTTGGAAACCAGATCATCCCTGA
CACTGCTATTCTGTCGGTGGTGCCCTTTCATCATGGGTTTGGGATGTT
CACAACACTGGGATACCTCATTTGCGGGTTTAGAGTGGTGCTCATGTA
TAGGTTTGAAGAAGAACTATTCCTACGCTCTTTGCAAGATTATAAGAT
15 TCAGTCTGCTCTGCTGGTGCCAACACTATTCTCTTTTTTTGCTAAGTCT
ACGCTCATAGACAAGTATGACTTGTCCAACCTGCACGAGATTGCTTCT
GGCGGAGCACCTCTGTCTAAGGAGGTAGGTGAGGCTGTGGCTAAGCG
CTTTCATCTGCCTGGTATCAGACAGGGGTACGGGCTAACAGAAACAA
CTTCTGCTATTCTGATTACACCAGAGGGCGATGACAAACCCGGGGCT
20 GTAGGGAAAGTGGTGCCCTTTTTTGAAGCCAAAGTAGTTGATCTTGAT
ACCGGTAAGACACTAGGGGTGAACCAGCGTGGTGAACCTGTGTGTGCG
GGGCCCTATGATTATGTCGGGGTACGTTAACAACCCCGAAGCTACAA
ATGCTCTCATAGACAAGGACGGGTGGCTTCATAGCGGCGACATTGCC
TACTGGGACGAGGATGAGCATTTCTTCATCGTGGACAGACTGAAGTC
25 GTTGATCAAATACAAGGGGTATCAAGTAGCTCCTGCCGAGCTTGAGT
CCATTCTGCTTCAACACCCCAATATCTTCGATGCTGGGGTGGCTGGGC
TGCCTGATGATGATGCTGGAGAGCTGCCTGCTGCTGTAGTAGTGCTTG
AGCATGGTAAGACAATGACAGAGAAGGAGATCGTGGATTATGTGGCT
TCACAAGTGACAACAGCTAAGAACTCCGAGGTGGCGTTGTGTTTGT
30 GGATGAGGTGCCTAAAGGGCTCACTGGCAAGCTGGATGCCAGAAAA
ATTCGAGAGATTCTCATTAAGGCTAAGAAGGGTGGAAAGATTGCTGT
GTAATAGTTCTAGA (SEQ ID NO:94)

hluc+ver2B8 has the following sequence

AAAGCCACCATGGAaGATGCcAAaAAcATTAAGAAGGGGCCTGCTCCc
 TTcTAcCCTCTtGAaGATGGGACtGcTGGcGAGCAaCTtCAcAAaGCTATGA
 5 AGcGgTATGCTCTtGTGCCaGGgACAATTGCgTTcACgGATGCTCAcATTG
 AaGTaGAcATcACATAcGCTGAGTATTTTGAGATGTCgGTGcGgCTGGCa
 GAaGCTATGAAGcGcTATGGGCTGAATACAAAcCATAGAATTGTaGTGT
 GcagTGAGAAcTCgtTGCAGTTcTTTATGCCcGTGCTGGGGGCTCTcTTcAT
 cGGGGTGGCTGTGGCTCCTGCTAAcGAcATcTAcAAcGAGcGAGAGCTgt
 10 TGAAcTCgATGGGGATcTCTCAGCCTACAGTGGTGTtTGTGagTAAGAA
 aGGGCTtCAaAAGATTCTcAATGTGCAaAAGAAGCTaCCgATcATaCAaAA
 GATcATcATcATGGAtagcAAGACcGAcTAcCAGGGGTTTCAGTCcATGTA
 cACATTTGTaACcTCTCATCTGCCTCCTGGcTTcAAcGAGTAcGAcTTcGT
 GCCcGAGTCTTTcGAcAGgGAcAAaACgATTGCTCTGATcATGAACagcTCc
 15 GGGTCTACcGGGCTGCCTAAGGGtGTaGCTCTGCCcCATcGAACAGCTT
 GTGTGAGATTcTCTCATGCcAGgGAcCCgATcTTtGGaAAcCAGATcATcC
 CTGAcAcTgCTATTCTGTCgGTgGTGCCcTTTCATCATGGGTTTGGGATG
 TTcACAACACTGGGaTAccTcATtTGcGGGTTTAGAGTGGTGCTcATGTAT
 AGgTTTGAaGAaGAaCTaTTccTacGcTCTtTGCAaGATTATAAGATTcAGT
 20 CTGCTCTGCTGGTGCCaACACTaTTcTCTTTTTTTGCTAAGTCTACgCTcA
 TaGAcAAGTATGActTGTCCaAAcTGCACgAGATTGCTTCTGGcGGaGCaCC
 TCTGTCTAAGGAGGTaGGtGAGGCTGTGGCTAAGcGcTTTCATCTGCCT
 GGtATcAGACAGGGGTAcGGGCTaACAGAAACAACtTCTGCTATTCTGAT
 TACACCgAGGGcGATGAcAAaCCtGGGGCTGTaGGGAaAGTGGTGCCcT
 25 TTTTGAaGCcAAaGTaGTtGATCTtGATACcGGtAAGACACTaGGGGTGA
 AcCAGcGtGGtGAaCTGTGTGTGcGgGGcCCTATGATTATGTCgGGGTAcG
 TtAAcAAcCCcGAaGCTACAAATGCTCTcATaGAcAAGGAcGGgTGGcTtC
 ATagcGGcGAcATTGCcTAcTGGGAcGAGGATGAGCATTtTTcATcGTGG
 AcAGACTGAAGTCgtTGATcAAaTAcAAGGGGTATCAaGTaGCTCCTGCc
 30 GAGCTtGAgTCcATTCTGCTtCAaCAcCCcAAAtATcTTcGATGCTGGGGTGG
 CTGGGCTGCCTGATGATGATGCTGGaGAGcTGCCTGCTGCTGTaGTaGT
 GCTtGAGCAtGGtAAGACAATGACAGAGAAGGAGATcGTGGATTATGT
 GGCTTCaCAaGTGACAACAGCTAAGAAaCTccGAGGtGGcGTtGTGTTTG
 TGGATGAGGTGCCTAAaGGGCTcActGGcAAGCTGGATGCcAGAAAaAT

TcGAGAGATTCTcATTAAGGCTAAGAAGGGtGGaAAGATTGCTGTGTA
ATAgTTCTAGA (SEQ ID NO:31).

hluc+ver2BF8 was modified to yield hluc+ver2BF9.

5

hluc+ver2B9 has the following sequence

AAAGCCACCATGGAaGATGCcAAaAAcATTAAGAAGGGGCCTGCTCCc
TTcTAcCCTCTtGAaGATGGGACtGCtGGcGAGCAaCTtCAcAAaGCTATGA
AGcGgTATGCTCTtGTGCCaGGgACAATTGCgTTcACgGATGCTCAcATTG
10 AaGTaGAcATcACATAcGCTGAGTATTTTGAGATGTCgGTGcGgCTGGCa
GAaGCTATGAAGcGcTATGGGCTGAATACAAAcCATAGAATTGTaGTGT
GcagTGAGAAcTCgtTGCAGTTcTTTATGCCcGTGCTGGGGGCTCTcTTcAT
tGGGGTGGCTGTGGCTCCTGCTAAtGAcATcTAcAAcGAGcGAGAGCTgtT
GAAcagtATGGGGATcTCTCAGCCTACAGTGGTGTtTGTAAGAAaG
15 GGCTtCAaAAGATTCTcAATGTGCAaAAGAAGCTaCCgATcATaCAaAAG
ATcATcATcATGGAtagcAAGACcGAcTAcCAGGGGTTTCAGTCcATGTAc
ACATTTGTaAcCTCTCATCTGCCTCCTGGcTTcAAtGAGTAAtGAcTTcGTG
CCcGAGTCTTTcGAcAGgGAcAAaACgATTGCTCTGATcATGAAcagcagtG
GGTCTACcGGGCTGCCTAAGGGtGTaGCTCTGCCcCATcGAACAGCTTG
20 TGTGAGATTcTCTCATGCcAGgGAcCCgATcTTtGGaAAcCAGATcATcCCT
GAcACtGCTATTCTGTcGtGtGTGCCcTTTCATCATGGGTTTGGGATGTT
cACAACACTGGGaTAccTcATtTGcGGGTTTAGAGTGGTGCTcATGTATA
GgTTTGAaGAaGAaCTaTTccTacGcTCTtTGCAaGATTATAAGATTcAGTC
TGCTCTGCTGGTGCCaACACTaTTcTCTTTTTTTGCTAAGTCTACgTcAT
25 aGAcAAGTATGActTGTCcAAcTGCAcGAGATTGCTTCTGGcGGaGCaCCT
CTGTCTAAGGAGGTaGGtGAGGCTGTGGCTAAGcGcTTTCATCTGCCTG
GtATcAGACAGGGGTAcGGGCTaACAGAAcACAACtTCTGCTATTCTGATT
ACACCaGAGGGcGATGAcAAaCCtGGGGCTGTaGGGAAaGTGGTGCCcTT
TTTTGAaGCcAAaGTaGtTgATCTtGATACcGGtAAGACACTaGGGGTGAA
30 cCAGaGaGGtGAatTGTGTGTGaGgGGcCCTATGATTATGTCgGGGTAcGtT
AAcAAcCCcGAaGCTACAAATGCTCTcATaGAcAAGGAcGGgTGGcTtCAT
agtGGaGAtATTGCcTAcTGGGATGAaGATGAGCATTtTcATcGTGGAcA
GACTGAAGTCgtTGATcAAaTAcAAGGGGTATCAaGTaGCTCCTGCcGAG

CTtGAgTCcATTCTGCTtCAaCAcCCcAAAtATcTTcGATGCTGGGGTGGCTG
 GGCTGCCTGATGATGATGCTGGaGAGcTGCCTGCTGCTGTaGTaGTGCTt
 GAGCAtGGtAAGACAATGACAGAGAAGGAGATcGTGGATTATGTGGCT
 TCaCAaGTGACAACAGCTAAGAAaCTccGAGGtGGcGTtGTGTTTGTGGA
 5 TGAGGTGCCTAAaGGGCTcActGGcAAGCTGGATGCcAGAAAaATTcGA
 GAGATTCTcATTAAGGCTAAGAAGGGtGGaAAGATTGCTGTGTAAATAgT
 TCTAGA (SEQ ID NO:32).

The *Bgl*I sequence in hluc+ver2BF9 was removed resulting in hluc+ver2BF10.

10 hluc+ver2BF10 demonstrated poor expression.

hluc+ver2B10 has the following sequence

AAAGCCACCATGGAaGATGCcAAaAAcATTAAGAAGGGGCCTGCTCCc
 TTcTAcCCTCTtGAaGATGGGACtGCTGGcGAGCAaCTtCAcAAaGCTATGA
 15 AGcGgTATGCTCTtGTGCCaGGgACAATTGCgTTcACgGATGCTCAcATTG
 AaGTaGAcATcACATAcGCTGAGTATTTTGAGATGTCgGTGcGgCTGGCa
 GAaGCTATGAAGcGcTATGGGCTGAATACAAAcCATAGAATTGTaGTGT
 GcagTGAGAAcTCgtTGCAGTTcTTTATGCCcGTGCTGGGGGCTCTcTTcAT
 tGGGGTGGCTGTGGCTCCTGCTAAAtGAcATcTAcAAcGAGcGAGAGCTgtT
 20 GAAcagtATGGGGATcTCTCAGCCTACAGTGGTGTtTGTGagTAAGAAaG
 GGCTtCAaAAGATTCTcAATGTGCAaAAGAAGCTaCCgATcATaCAaAAG
 ATcATcATcATGGAtagcAAGACcGAcTAcCAGGGGTTTCAGTCcATGTAc
 ACATTTGTaACcTCTCATCTGCCTCCTGGcTTcAAAtGAGTAAtGAcTTcGTG
 CCcGAGTCTTTcGAcAGgGAcAAaACgATTGCTCTGATcATGAAcagcagtG
 25 GGTCTACcGGGCTGCCTAAGGGtGTaGCTCTGCCcCATcGAACAGCTTG
 TGTGAGATTcTCTCATGCcAGgGAcCCgATcTTtGGaAAcCAGATcATcCCT
 GAcAcTgCTATTCTGTCgGTgGTGCCcTTTCATCATGGGTTTGGGATGTT
 cACAACACTGGGaTAccTcATtTgcGGGTTTAGAGTGGTGCTcATGTATA
 GgTTTGAaGAaGAaCTaTTccTacGcTCTtTGCAaGATTATAAGATTcAGTC
 30 TGCTCTGCTGGTGCCaACACTaTTcTCTTTTTTTGCTAAGTCTACgCTcAT
 aGAcAAGTATGActGTTCcAAcTGCAcGAGATTGCTTCTGGcGGaGCcCCT
 CTGTCTAAGGAGGTaGGtGAGGCTGTGGCTAAGcGcTTTCATCTGCCTG
 GtATcAGACAGGGGTAcGGGCTaACAGAAACAACtCTGCTATTCTGATT

ACACCaGAGGGcGATGAcAAaCCtGGGGCTGTaGGGAAaGTGGTGCCcTT
 TTTTGAAgCcAAaGTaGTtGATCTtGATACcGGtAAGACACTaGGGGTGAA
 cCAGaGaGGtGAatTGTGTGTGaGgGGcCCTATGATTATGTCgGGGTAcGTt
 AAcAAcCCcGAaGCTACAAATGCTCTcATaGAcAAGGAcGGgTGGcTtCAT
 5 agtGGaGAtATTGCcTAcTGGGAtGAaGATGAGCATTtTtCATcGTGGAcA
 GACTGAAGTCgtTGATcAAaTAcAAGGGGTATCAaGTaGCTCCTGCcGAG
 CTtGAgTCcATTCTGCTtCAaCAcCCcAAaTtTcGATGCTGGGGTGGCTG
 GGCTGCCTGATGATGATGCTGGaGAGcTGCCTGCTGCTGTaGTaGTGCTt
 GAGCAtGGtAAGACAATGACAGAGAAGGAGATcGTGGATTATGTGGCT
 10 TCaCAaGTGACAACAGCTAAGAAaCTccGAGGtGGcGTtGTGTTTGTGGA
 TGAGGTGCCTAAaGGaCTcACtGGcAAGCTGGATGCcAGAAAaATTcGAG
 AGATTCTcATTAAGGCTAAGAAGGGtGGaAAGATTGCTGTGTAATAgTT
 CTAGA (SEQ ID NO:33).

15

Table 11

Summary of Firefly Luciferase Constructs

| Firefly luciferase Gene | Number of consensus transcription factor binding sites | Number of Promoter modules* | CG dinucleotides (possible methylation sites) |
|-------------------------|--|-----------------------------|---|
| Luc+ | 287 | 7 | 97 |
| hluc+ver2AF8 | 3 | 0 | 132 |
| hluc+ver2BF10 | 3 | 0 | 43 |

*Promoter modules are defined as a composite regulatory element, with 2 TFBS separated by a spacer, which has been shown to exhibit synergistic or
 20 antagonistic function.

Example 4Synthetic Selectable Polypeptide GenesDesign Process25 Define sequences

Protein sequence that should be maintained:

- Neo: from *neo* gene of pCI-neo (Promega) (SEQ ID NO:1)
- Hyg: from *hyg* gene of pcDNA3.1/Hygro (Invitrogen) (SEQ ID NO:6)

DNA flanking regions for starting sequence:

- 5' end: Kozak sequence from *neo* gene of pCI-neo (GCCACCATGA; SEQ ID NO:34)), *Pf*MI site (CCANNNNTGG; SEQ ID NO:35), add Ns at end (to avoid search algorithm errors & keep ORF1):
neo/hyg: NNNNNCCAnnnnnTGGCCACC-ATG-G (SEQ ID NO:36)
- 5 Change: replace *Pf*MI with *Sbf*I (CCTGCAGG)
- 3' end: two stop codons (at least one TAA), *Pf*MI site (not compatible with that at 5' end to allow directional cloning), add Ns at end (to avoid search algorithm errors):
neo/hyg: TAATAACCAnnnnnTGGNNN (SEQ ID NO:37)
- 10 Change: replace *Pf*MI with *Afl*III (CTTAAG)

Define codon usage

Codon usage was obtained from the Codon Usage Database

(<http://www.kazusa.or.jp/codon/>):

- 15 Based on: GenBank Release 131.0 [15 August 2002] (Nakamura et al., 2000).

Codon usage tables were downloaded for:

- HS: *Homo sapiens* [gbpri] 50,031 CDS's (21,930,294 codons)
- MM: *Mus musculus* [gbrod] 23,113 CDS's (10,345,401 codons)
- 20 EC: *Escherichia coli* [gbtct] 11,985 CDS's (3,688,954 codons)
- EC K12: *Escherichia coli* K12 [gbtct] 4,291 CDS's (1,363,716 codons)
- ⇒ HS and MM were compared and found to be closely similar, use HS table
- ⇒ EC and EC K12 were compared and found to be closely similar, use
- 25 EC K12 table

Codon selection strategy:

- Overall strategy is to adapt codon usage for optimal expression in mammalian cells while avoiding low-usage *E. coli* codons. One "best" codon was selected for each amino acid and used to back-translate the
- 30 desired protein sequence to yield a starting gene sequence.

Strategy A was chosen for the design of the *neo* and *hyg* genes (see Table 12). (Strategy A: Codon bias optimized: emphasis on codons showing the highest usage frequency in HS. Best codons are those with highest

usage in HS, unless a codon with slightly lower usage has substantially higher usage in *E. coli*).

Table 12

| Amino acid | Codon Choices in Examples 1-2 | Codon Choices in Codon Bias Optimized Strategy A |
|------------|-------------------------------|--|
| Gly | GGC/GGT | GGC |
| Glu | GAG | GAG |
| Asp | GAC | GAC |
| Val | GTG/GTC | GTG |
| Ala | GCC/GCT | GCC |
| Arg | CGC/CGT | CGC |
| Ser | TCT/AGC | AGC |
| Lys | AAG | AAG |
| Asn | AAC | AAC |
| Ile | ATC/ATT | ATC |
| Thr | ACC/ACT | ACC |
| Cys | TGC | TGC |
| Tyr | TAC | TAC |
| Leu | CTG/TTG | CTG |
| Phe | TTC | TTC |
| Gln | CAG | CAG |
| His | CAC | CAC |
| Pro | CCA/CCT | CCC |

5

Generate starting gene sequences

Use custom codon usage table in Vector NTI 8.0 (Informax) ("Strategy A")

Back-translate *neo* and *hyg* protein sequences

Neo (based on neomycin gene from Promega's pCI-neo)

10 MIEQDGLHAGSPAAWVERLFGYDWAQQTIGCSDAAVFRLSAQGRPVL
 VKTDLSGALNELQDEAARLSWLATTGVPCAAVLDVVTEAGRDWLLGE
 VPGQDLLSSHLAPAEKVSIMADAMRRLHTLDPATCPFDHQAKHRIERAR

TRMEAGLVDQDDLDEEHQGLAPAE L FARLKARMPDGEDLVVTHGDAC
LPNIMVENGRFSGFIDCGRLGVADRYQDIALATRDIAEELGGEWADRFLV
LYGIAAPDSQRIAFYRLLDEFF (SEQ ID NO:2) and encoded by

Atgattgaacaagatggattgcacgcaggttctccggccgcttgggtggagaggctattcggctatgactgggcac
5 aacagacaatcggtgctctgatgccgctgttccggctgtcagcgcaggggcgcccggttctttgtcaagacc
gacctgtccggtgccctgaatgaactgcaggacgaggcagcgcggctatcgtggctggccacgacgggcgttcct
tgcgcagctgtctgcaggtgtcactgaagcgggaaggactggctgctattgggcgaagtgcggggcaggat
ctcctgtcatctcacctgtctctgccgagaaagtatccatcatggctgatgcaatgcggcggtgcatacgcttgatc
cggctacctgcccattcgaccaccaagcgaacatcgcatcgagcgagcacgtactcggatggaagccggtcttgt
10 cgatcaggatgatctggacgaagagcatcaggggctcgcgccagccgaactgttcgccaggctcaaggcgcgcgat
gcccgcggcgaggatctcgtcgtgacctatggcgatgcctgcttgcgaatatcatggtgaaaatggccgctttt
ctggattcatcgactgtggccggctgggtgtggcggaccgctatcaggacatagcgttggtacctcgatattgctg
aagagcttggcggcgaatgggctgaccgcttctcgtgctttacggatcgccgctcccattcgcagcgcatcgcc
ttctatcgcttcttgacgagttcttctga (SEQ ID NO:1)

15 Hyg (based on hygromycin gene from Invitrogen's pcDNA3.1/Hygro)
MKKPELTATSVEKFLIEKFD SVSDLMQLSEGEESRAFSFDVGGRGYVLRV
NSCADGFYKDRYVYRHFAS AALPIPEVLDIGEFSESLTYCISRR AQGVTLQ
DLPETELPAVLQPVAEAMD AIAAADLSQTSFGFPGFPQGIGQYTTWRDFI
20 CAIADPHVYHWQTVMDDTV SASVAQALDELMLWAEDCPEVRHLVHAD
FGSNNVLTDNGRITAVIDWSEAMFGDSQYEVANIFFWRPWLACMEQQT
RYFERRHPELAGSPRLRAYMLRIGLDQLYQSLVDGNFDDAAWAQGRCD
AIVRSGAGTVGRTQIARRSA AVWTDGCVEVLADSGNRRPSTRPRAKE
(SEQ ID NO:7) encoded by

25 Atgaaaaagcctgaactcaccgcgacgtctgtcgagaagttctgatcgaagttcgacagcgtctccgacctgat
gcagctctcgagggcgaagaatctcgtgctttcagctcgtatgtaggagggcgtggatatgtcctgcgggtaaata
gctgcgccgatggtttctacaaagatcgttatgtttatcggcactttgcatcgccgcgctcccattccggaagtgtt
gacattggggaattcagcgagagcctgacctattgcatctcccgctgcacagggtgtcacgttgcaagacctgcc
30 tgaaccggaactgcccgtgttctgcagccggctcgcggaggccatggatgcgatcgctgcggccgatcttagccag
acgagcgggttcggccattcgaccgcaaggaatcggtcaatacactacatggcgtgattcatatgcgcgattgc
tgatccccatgtgtacttggaactgtgatggacgacaccgtcagtcgtccgtcgcgcaggctctgatgagc
tgatgctttgggcccaggactgccccgaagtccggcacctcgtgcacgcggatttcgggtccaacaatgtcctgacg

gacaatggccgcataacagcgggtcattgactggagcgaggcgatgttcggggattcccaatacagaggtcgccaac
 atcttctctggaggccgtgggttgctgtatggagcagcagacgcgctacttcgagcggaggcatccggagcttgc
 aggatcgccgcgggtccggcgatatgctccgcatgtgttgacctatcagagcttgggtgacggcaattc
 gatgatgcagcttgggcgcagggtcgatgcgacgcaatgtccgatccggagccgggactgtcggcgtagacaaa
 5 atcggccgcagaagcgcggccgtctggaccgatggctgtgtagaagtactcgccgatagtggaaaccgacgcccc
 agcactcgtccgagggcaaaggaat (SEQ ID NO:6).

Table 13

Nomenclature of exemplary neo and hyg gene versions

| Gene name | Description |
|-----------|--|
| neo | from pCI-neo (Promega) |
| hneo | humanized (codon usage strategy A) ORF |
| hneo-F | humanized ORF with 5' and 3' flanking regions |
| hneo-1F | humanized ORF with 5' and 3' flanking regions after first removal of undesired sequence matches |
| hneo-2F | humanized ORF with 5' and 3' flanking regions after second removal of undesired sequence matches |
| hneo-3F | humanized ORF with 5' and 3' flanking regions after third removal of undesired sequence matches |
| hneo-3FB | Changed 5' and 3' flanking cloning sites |
| hyg | from pcDNA3.1/Hygro (Invitrogen) |
| hhyg | humanized (codon usage strategy A) ORF |
| hhyg-F | humanized ORF with 5' and 3' flanking regions |
| hhyg-1F | humanized ORF with 5' and 3' flanking regions after first removal of undesired sequence matches |
| hhyg-2F | humanized ORF with 5' and 3' flanking regions after second removal of undesired sequence matches |
| hhyg-3F | humanized ORF with 5' and 3' flanking regions after third removal of undesired sequence matches |
| hhyg-3FB | Changed 5' and 3' flanking cloning sites |

“h” indicates humanized codons, “F” indicates presence of 5' and 3' flanking sequences.

Create starting (codon-optimized) gene sequences:

- 5 hneo (humanized starting gene sequence without flanking regions in hneo-F)
 CCACTCAGTGGCCACCATGATCGAGCAGGACGGCCTGCACGCCGGCA
 GCCCCGCCGCTGGGTGGAGCGCCTGTTCCGGCTACGACTGGGCCAG
 CAGACCATCGGCTGCAGCGACGCCGCCGTGTTCCGCCTGAGCGCCCA
 GGGCCGCCCCGTGCTGTTCTGTGAAGACCGACCTGAGCGGCGCCCTGA
 10 ACGAGCTGCAGGACGAGGCCGCCGCTGAGCTGGCTGGCCACCACC
 GGCGTGCCCTGCGCCGCCGTGCTGGACGTGGTGACCGAGGCCGGCCG
 CGACTGGCTGCTGCTGGGCGAGGTGCCCGGCCAGGACCTGCTGAGCA
 GCCACCTGGCCCCCGCCGAGAAGGTGAGCATCATGGCCGACGCCATG
 CGCCGCTGCACACCCTGGACCCCGCCACCTGCCCCTTCGACCACCA
 15 GGCCAAGCACCGCATCGAGCGCGCCCGCACCCGCATGGAGGCCGGC
 CTGGTGACCAAGGACGACCTGGACGAGGAGCACCAGGGCCTGGCCC
 CCGCCGAGCTGTTCCGCCGCTGAAGGCCCGCATGCCCGACGGCGAG
 GACCTGGTGCTGACCCACGGCGACGCCTGCCTGCCCAACATCATGGT
 GGAGAACGGCCGCTTCAGCGGCTTCATCGACTGCGGCCGCTGGGCG
 20 TGGCCGACCGCTACCAGGACATCGCCCTGGCCACCCGCGACATCGCC
 GAGGAGCTGGGCGGCGAGTGGGCCGACCGCTTCCTGGTGCTGTACGG
 CATCGCCGCCCCCGACAGCCAGCGCATCGCCTTCTACCGCCTGCTGG
 ACGAGTTCTTCTAATAACCAGTCTCTGG (SEQ ID NO:3).
- 25 hhyg (humanized starting gene sequence without flanking regions)
 CCACTCAGTGGCCACCATGAAGAAGCCCGAGCTGACCGC CACCAGCG
 TGGAGAAGTTCCTGATCGAGAAGTTCGACAGCGTGAGCGACCTGATG
 CAGCTGAGCGAGGGCGAGGAGAGCCGCGCCTTCAGCTTCGACGTGG
 GCGGCCGCGGCTACGTGCTGCGCGTGAACAGCTGCGCCGACGGCTTC
 30 TACAAGGACCGCTACGTGTACCGCCACTTCGCCAGCGCCGCCCTGCC
 CATCCCCGAGGTGCTGGACATCGGCGAGTTCAGCGAGAGCCTGACCT
 ACTGCATCAGCCGCCGCGCCCAGGGCGTGACCCTGCAGGACCTGCCC
 GAGACCGAGCTGCCCCGCCGTGCTGCAGCCCGTGGCCGAGGCCATGGA

CGCCATCGCCGCCGCCGACCTGAGCCAGACCAGCGGCTTCGGCCCCCT
TCGGCCCCCAGGGCATCGGCCAGTACACCACCTGGCGCGACTTCATC
TGCGCCATCGCCGACCCCCACGTGTACCACTGGCAGACCGTGATGGA
CGACACCGTGAGCGCCAGCGTGGCCCAGGCCCTGGACGAGCTGATGC
5 TGTGGGCCGAGGACTGCCCCGAGGTGCGCCACCTGGTGACGCCGAC
TTCGGCAGCAACAACGTGCTGACCGACAACGGCCGCATCACCGCCGT
GATCGACTGGAGCGAGGCCATGTTTCGGCGACAGCCAGTACGAGGTGG
CCAACATCTTCTTCTGGCGCCCCCTGGCTGGCCTGCATGGAGCAGCAG
ACCCGCTACTTCGAGCGCCGCCACCCCGAGCTGGCCGGCAGCCCCCG
10 CCTGCGCGCCTACATGCTGCGCATCGGCCTGGACCAGCTGTACCAGA
GCCTGGTGGACGGCAACTTCGACGACGCCGCCTGGGCCCAGGGCCGC
TGCGACGCCATCGTGCGCAGCGGCCCGGCACCGTGGGCCGCACCCA
GATCGCCCGCCGACGCGCCGCGTGTGGACCGACGGCTGCGTGGAGG
TGCTGGCCGACAGCGGCAACCGCCGCCCCAGCACCCGCCCCCGCGCC
15 AAGGAGTAATAACCAGCTCTTGG (SEQ ID NO:8).

Programs and databases used for identification and removal of sequence motifs

All from Genomatix Software GmbH (Munich, Germany,

<http://www.genomatix.de>):

GEMS Launcher Release 3.5.2 (June 2003)

- 20 MatInspector professional Release 6.2.1 June 2003
Matrix Family Library Ver 3.1.2 June 2003 (incl. 318 vertebrate matrices
in 128 families)
ModelInspector professional Release 4.8 October 2002
Model Library Ver 3.1 March 2003 (226 modules)
25 SequenceShaper tool
User Defined Matrices

Sequence motifs to remove from starting gene sequences

(In order of priority)

- 30 Restriction enzyme recognition sequences:
See user-defined matrix subset *neo* and *hyg*. Same as those used for
design of hluc+ version 2.0
Generally includes those required for cloning (pGL4) or commonly used

for cloning

Change: also *Sbf*I, *Afl*II, *Acc*III

Transcription factor binding sequences:

Promoter modules (2 TF binding sites with defined orientation) with

5 default score or greater

Vertebrate TF binding sequences with score of at least core=0.75 /

matrix=optimized

Eukaryotic transcription regulatory sites:

Kozak sequence

10 Splice donor / acceptor sequences in (+) strand

PolyA addition sequences in (+) strand

Prokaryotic transcription regulatory sequences:

E. coli promoters

E. coli RBS (if less than 20 bp upstream of Met codon)

15

User-defined matrix subset “*neo+hyg*”

Format: Matrix name (core similarity threshold / matrix similarity threshold)

- U\$AatII (0.75/1.00)
- U\$BamHI (0.75/1.00)
- 20 • U\$BglII (0.75/1.00)
- U\$BglIII (0.75/1.00)
- U\$BsaI (0.75/1.00)
- U\$BsmAI (0.75/1.00)
- U\$BsmBI (0.75/1.00)
- 25 • U\$BstEII (0.75/1.00)
- U\$BstXI (0.75/1.00)
- U\$Csp45I (0.75/1.00)
- U\$CspI (0.75/1.00)

- U\$EC-P-10 (1.00/Optimized)
- U\$EC-P-35 (1.00/Optimized)
- U\$EC-Prom (1.00/Optimized)
- U\$EC-RBS (0.75/1.00)
- 5 • U\$EcoRI (0.75/1.00)
- U\$HindIII (0.75/1.00)
- U\$Kozak (0.75/Optimized)
- U\$KpnI (0.75/1.00)
- U\$MluI (0.75/1.00)
- 10 • U\$NcoI (0.75/1.00)
- U\$NdeI (0.75/1.00)
- U\$NheI (0.75/1.00)
- U\$NotI (0.75/1.00)
- U\$NsiI (0.75/1.00)
- 15 • U\$PflMI (0.75/1.00)
- U\$PmeI (0.75/1.00)
- U\$PolyAsig (0.75/1.00)
- U\$PstI (0.75/1.00)
- U\$SacI (0.75/1.00)
- 20 • U\$SacII (0.75/1.00)
- U\$SalI (0.75/1.00)
- U\$SfiI (0.75/1.00)
- U\$SgfI (0.75/1.00)

- U\$SmaI (0.75/1.00)
- U\$SnaBI (0.75/1.00)
- U\$SpeI (0.75/1.00)
- U\$Splice-A (0.75/Optimized)
- 5 • U\$Splice-D (0.75/Optimized)
- U\$XbaI (0.75/1.00)
- U\$XcmI (0.75/1.00)
- U\$XhoI (0.75/1.00)
- ALL vertebrates.lib (0.75/Optimized)

10

User-defined matrix subset “*neo+hyg-EC*”

Format: Matrix name (core similarity threshold / matrix similarity threshold)

- U\$AatII (0.75/1.00)
- U\$BamHI (0.75/1.00)
- 15 • U\$BglII (0.75/1.00)
- U\$BglIII (0.75/1.00)
- U\$BsaI (0.75/1.00)
- U\$BsmAI (0.75/1.00)
- U\$BsmBI (0.75/1.00)
- 20 • U\$BstEII (0.75/1.00)
- U\$BstXI (0.75/1.00)
- U\$Csp45I (0.75/1.00)
- U\$CspI (0.75/1.00)
- U\$EcoRI (0.75/1.00)

- U\$HindIII (0.75/1.00)
- U\$Kozak (0.75/Optimized)
- U\$KpnI (0.75/1.00)
- U\$MluI (0.75/1.00)
- 5 • U\$NcoI (0.75/1.00)
- U\$NdeI (0.75/1.00)
- U\$NheI (0.75/1.00)
- U\$NotI (0.75/1.00)
- U\$NsiI (0.75/1.00)
- 10 • U\$PflMI (0.75/1.00)
- U\$PmeI (0.75/1.00)
- U\$PolyAsig (0.75/1.00)
- U\$PstI (0.75/1.00)
- U\$SacI (0.75/1.00)
- 15 • U\$SacII (0.75/1.00)
- U\$SalI (0.75/1.00)
- U\$SfiI (0.75/1.00)
- U\$SgfI (0.75/1.00)
- U\$SmaI (0.75/1.00)
- 20 • U\$SnaBI (0.75/1.00)
- U\$SpeI (0.75/1.00)
- U\$Splice-A (0.75/Optimized)
- U\$Splice-D (0.75/Optimized)

- U\$XbaI (0.75/1.00)
- U\$XcmI (0.75/1.00)
- U\$XhoI (0.75/1.00)
- ALL vertebrates.lib (0.75/Optimized)

5

User-defined matrix subset "pGL4-072503"

Format: Matrix name (core similarity threshold / matrix similarity threshold)

- U\$AatII (0.75/1.00)
- U\$AccIII (0.75/1.00)
- 10 • U\$AflIII (0.75/1.00)
- U\$BamHI (0.75/1.00)
- U\$BglII (0.75/1.00)
- U\$BglIII (0.75/1.00)
- U\$BsaI (0.75/1.00)
- 15 • U\$BsmAI (0.75/1.00)
- U\$BsmBI (0.75/1.00)
- U\$BstEII (0.75/1.00)
- U\$BstXI (0.75/1.00)
- U\$Csp45I (0.75/1.00)
- 20 • U\$CspI (0.75/1.00)
- U\$EC-P-10 (1.00/Optimized)
- U\$EC-P-35 (1.00/Optimized)
- U\$EC-Prom (1.00/Optimized)
- U\$EC-RBS (0.75/1.00)

- U\$EcoRI (0.75/1.00)
- U\$HindIII (0.75/1.00)
- U\$Kozak (0.75/Optimized)
- U\$KpnI (0.75/1.00)
- 5 • U\$MluI (0.75/1.00)
- U\$NcoI (0.75/1.00)
- U\$NdeI (0.75/1.00)
- U\$NheI (0.75/1.00)
- U\$NotI (0.75/1.00)
- 10 • U\$NsiI (0.75/1.00)
- U\$PflMI (0.75/1.00)
- U\$PmeI (0.75/1.00)
- U\$PolyAsig (0.75/1.00)
- U\$PstI (0.75/1.00)
- 15 • U\$SacI (0.75/1.00)
- U\$SacII (0.75/1.00)
- U\$SalI (0.75/1.00)
- U\$SbfI (0.75/1.00)
- U\$SfiI (0.75/1.00)
- 20 • U\$SgfI (0.75/1.00)
- U\$SmaI (0.75/1.00)
- U\$SnaBI (0.75/1.00)
- U\$SpeI (0.75/1.00)

- U\$Splice-A (0.75/Optimized)
- U\$Splice-D (0.75/Optimized)
- U\$XbaI (0.75/1.00)
- U\$XcmI (0.75/1.00)
- 5 • U\$XhoI (0.75/1.00)
- ALL vertebrates.lib

Strategy for removal of sequence motifs

The undesired sequence motifs specified above were removed from the
 10 starting gene sequence by selecting alternate codons that allowed retention of the
 specified protein and flanking sequences. Alternate codons were selected in a
 way to conform to the overall codon selection strategy as much as possible.

General steps:

- 15 - Identify undesired sequence matches with MatInspector using matrix family
 subset "neo+hyg" or "neo+hyg-EC" and with ModelInspector using default
 settings.
- Identify possible replacement codons to remove undesired sequence matches
 with SequenceShaper (keep ORF).
- 20 - Incorporate changes into a new version of the synthetic gene sequence and
 re-analyze with MatInspector and ModelInspector.

Specific steps:

- First try to remove undesired sequence matches using subset "neo+hyg-EC"
 and SequenceShaper default remaining thresholds (0.70/Opt-0.20).
- 25 - For sequence matches that cannot be removed with this approach use lower
 SequenceShaper remaining thresholds (e.g. 0.70/Opt-0.05).
- For sequence matches that still cannot be removed, try different
 combinations of manually chosen replacement codons (especially if more
 than 3 base changes might be needed). If that introduces new sequence

matches, try to remove those using the steps above (a different starting sequence sometimes allows a different removal solution).

- Use subset "neo+hyg" to check whether problematic *E. coli* sequence matches were introduced, and if so try to remove them using an analogous approach to that described above for non *E. coli* sequences.

Use an analogous strategy for the flanking (non-ORF) sequences.

Final check with subset "pGL4-072503" after change in flanking cloning sites

After codon optimizing *neo* and *hyg*, hneo and hhyg were obtained.

- 10 Regulatory sequences were removed from hneo and hhyg yielding hneo-1F and hhyg-1F (the corresponding sequences without flanking regions are SEQ ID Nos. 38 and 30, respectively). Regulatory sequences were removed from hneo-1F and hhyg-1F yielding hneo-2F and hhyg-2F (the corresponding sequences without flanking regions are SEQ ID Nos. 39 and 42, respectively). Regulatory sequences were removed from hneo-2F and hhyg-2F yielding hneo-3F and hhyg-3F. Hneo-3F and hhyg-3F were further modified by altering 5' and 3' cloning sites yielding hneo-3FB and hhyg-3FB:

- hneo-3 (after 3rd round of sequence removal, subset neo+hyg) has the following sequence:

20 CCACTCcGTGGCCACCA TGATCGAaCAaGACGGCCTcCAtGCtGGCAGtC
CCGCaGCtTGGGTcGAaCGCtGTTCGGgTACGACTGGGCCCAGCAGAC
CATCGGaTGtAGCGAtGCgGCCGTGTTCCGtCTaAGCGCtCAaGGCCGgCC
CGTGCTGTTCTGAAGACCGACCTGAGCGGCGCCCTGAACGAGCTtCA
25 aGACGAGGctGCCCCGCTGAGCTGGCTGGCCACCACCGGtGTaCCCTGC
GCCGctGTgtTGGAtGTtGTGACCGAaGCCGGCCGgGACTGGCTGCTGCT
GGGCGAGGTcCCtGGCCAGGAAtCTGCTGAGCAGCCACCTtGCCCCGct
GAGAAGGTtCtCATCATGGCCGAtGCaATGCGgCGCCTGCACACCCTGG
ACCCCGctACaTGCCCC TTCGACCACCAGGctAAGCAtCGgATCGAGCGt
30 GCtCGgACCCGCATGGAAGGCCGGCCTGGTGGACCAGGACGACCTGGA
CGAGGAGCAtCAGGGCCTGGCCCCCGctGAaCTGTTCGCCCCGCTGAAa
GCCCGCATGCCgGACGgtGAGGACCTGGTtGTGACaCAtGGtGAtGCCTG
CCTcCCtAACATCATGGTcGAGAAAtGGcCGCTTctCGGCTTCATCGACTG

CGGtCGCCTaGGaGtGCCGACCGCTACCAGGACATCGCCCTGGCCACC
 CGCGACATCGCtGAGGAGCTtGGCGGCGAGTGGGCCGACCGCTTcTaG
 TctTGTAACGGCATCGCaGcCCCGACAGCCAGCGCATCGCCTTCTACCG
 CCTGCTcGACGAGTTCTTtTAATGACCAGgCTCTGG (SEQ ID NO:4);

- 5 hneo-3FB (change *Pfl*MI sites to *Sbf*I at 5' end and *Afl*III at 3' end) has the following sequence:

cctgcaggCCACCATGATCGAACAAGACGGCCTCCATGCTGGCAGTCCCG
 CAGCTTGGGTCTGAACGCTTGTTCTGGGTACGACTGGGCCCAGCAGACC
 ATCGGATGTAGCGATGCGGCCGTGTTCCGTCTAAGCGCTCAAGGCCG
 10 GCCCGTGCTGTTCTGTGAAGACCGACCTGAGCGGCGCCCTGAACGAGC
 TTCAAGACGAGGCTGCCCCGCTGAGCTGGCTGGCCACCACCGGTGTA
 CCCTGCGCCGCTGTGTTGGA TGTTGTGACCGAAGCCGGCCGGGACTG
 GCTGCTGCTGGGCGAGGTCCCTGGCCAGGATCTGCTGAGCAGCCACC
 TTGCCCCCGCTGAGAAGGTTTCCATCATGGCCGATGCAATGCGGCGC
 15 CTGCACACCCTGGACCCCGCTACATGCCCCCTTCGACCACCAGGCTAA
 GCATCGGATCGAGCGTGCTCGGACCCGCATGGAGGCCGGCCTGGTGG
 ACCAGGACGACCTGGACGA GGAGCATCAGGGCCTGGCCCCCGCTGA
 ACTGTTCTGCCCCGCTGAAAGCCCCGCATGCCGGACGGTGAGGACCTGG
 TTGTGACACATGGTGATGCCTGCCTCCCTAACATCATGGTTCGAGAAT
 20 GGCCGCTTCTCCGGCTTCATCGACTGCGGTGCGCTAGGAGTTGCCGAC
 CGCTACCAGGACATCGCCCTGGCCACCCGCGACATCGCTGAGGAGCT
 TGGCGGCGAGTGGGCCGACCGCTTCTTAGTCTTGTACGGCATCGCAG
 CTCCCGACAGCCAGCGCATCGCCTTCTACCGCCTGCTCGACGAGTTCT
 TTTAATGAgcttaag (SEQ ID NO:5);

- 25 hhyg-3 (after 3rd round of sequence removal, subset neo+hyg) has the following sequence:

CCACTCcGTGGCCACCATGAAGAAGCCCGAGCTGACCGCtACCAGCGT
 tGAaAAaTTtCTcATCGAGAAGTTCGACAGtGTGAGCGACCTGATGCAGt
 TgtcgAGGGCGAaGAgAGCCGaGCCTTCAGCTTCGAtGTcGGCGGaCGC
 30 GGCTAtGTaCTGCGgGTGAAtAGCTGCGCtGAtGGCTTCTACAAaGACCG
 CTACGTGTACCGCCACTTCGCCAGCGCtGCaCTaCCCATCCCCGAaGTGt
 TGGACATCGGCGAGTTCAGCGAGAGCCTGACaTACTGCATCAGtaGaCG

CGCCCAaGGCGTtACtCTcCAaGAACCTcCCCGAaACaGAGCTGCCtGcTGT
 GtTaCAGCCTGTcGCCGAaGcTATGGAAtGcTATtGCCGCCGCCGACCTcAGt
 CAaACCAGCGGCTTCGGCCCaTTCGGgCCCCAaGGCATCGGCCAGTAC
 ACaACCTGGCGgGAAtTTCATtTGCGCCATtGcTGAAtCCCCAtGTcTACCACT
 5 GGCAGACCGTGATGGACGACA CCGTGtcCGCCAGCGTaGcTCAaGCCCT
 GGACGAaCTGATGCTGTGGGCCGAaGACTGtCCCGAGGTGCGCCAeCTc
 GTcCAtGCCGACTTCGGCAGCAACAACGTcCTGACCGACAACGGCCGC
 ATCACCGCCGTaATCGACTGGtcCGAaGcTATGTTCGGgGACAGtCAGTA
 CGAGGTGGCCAACATCTTCTTCTGGCGgCCCTGGCTGGCtTGCATGGA
 10 GCAGCAGACtCGCTACTTCGAGCGCCGgCAtCCCGAGCTGGCCGGCAG
 CCCtCGtCTGCGaGCCTACATGCTGCGCATCGGCCTGGAAtCAGCTcTACC
 AGAGCCTcGTGGACGGCAACTTCGACGAAtGcTGCCTGGGCtCAaGGCCG
 CTGCGAtGCCATCGTcCGCAGCGGgGCCGGCACCGTcGGtCGCACaCaAa
 TCGCtCGCCGgAGCGCCGCCGTaTGGACCGACGGCTGCGTcGAGGTGCT
 15 GGCCGACAGCGGCAACCGCCGgCCCAGtACaCGaCCgCGCGCtAAGGAG
 TAgTAACCAAggtcTGG (SEQ ID NO:9); and

hhyg-3FB (change *Pf*MI sites to *Sbf*I at 5' end and *Afl*II at 3' end) has the
 following sequence:

cctgcaggCCACCATGAAGAAGCCCGAGCTGACCGCTACCAGCGTTGAAA
 20 AATTTCTCATCGAGAAGTTCGACAGTGTGAGCGACCTGATGCAGTTG
 TCGGAGGGCGAAGAGAGCCGAGCCTTCAGCTTCGATGTGCGCGGACG
 CGGCTATGTACTGCGGGTGAA TAGCTGCGCTGATGGCTTCTACAAAG
 ACCGCTACGTGTACCGCCACTTCGCCAGCGCTGCACTACCCATCCCC
 GAAGTGTGGACATCGGCGAGTTCAGCGAGAGCCTGACATACTGCAT
 25 CAGTAGACGCGCCCAAGGCGTTACTCTCCAAGACCTCCCCGAAACAG
 AGCTGCCTGCTGTGTTACAGCCTGTGCGCGAAGCTATGGATGCTATTG
 CCGCCGCCGACCTCAGTCAAA CCAGCGGCTTCGGCCCATTCGGGCCC
 CAAGGCATCGGCCAGTACACAACCTGGCGGGATTTCATTTGCGCCAT
 TGCTGATCCCCATGTCTACCACTGGCAGACCGTGATGGACGACACCG
 30 TGTCCGCCAGCGTAGCTCAAGCCCTGGACGAACTGATGCTGTGGGCC
 GAAGACTGTCCCGAGGTGCGCCACCTCGTCCATGCCGACTTCGGCAG
 CAACAACGTCTCTGACCGACAAACGGCCGCATCACCGCCGTAATCGACT

GGTCCGAAGCTATGTTTCGGGGA CAGTCAGTACGAGGTGGCCAACATC
 TTCTTCTGGCGGCCCTGGCTGGCTTGCATGGAGCAGCAGACTCGCTAC
 TTCGAGCGCCGGCATCCCGAGCTGGCCGGCAGCCCTCGTCTGCGAGC
 CTACATGCTGCGCATCGGCCTGGATCAGCTCTACCAGAGCCTCGTGG
 5 ACGGCAACTTCGACGATGCTGC CTGGGCTCAAGGCCGCTGCGATGCC
 ATCGTCCGCAGCGGGGCCGGCA CCGTCGGTCGCACACAAATCGCTCG
 CCGGAGCGCCGCCGTATGGACC GACGGCTGCGTCGAGGTGCTGGCCG
 ACAGCGGCAACCGCCGGCCAGTACACGACCGCGCGCTAAGGAGTA
 GTAAActtaag (SEQ ID NO:10).

10 Analysis of hneo-3FB and hhyg-3FB

hneo-3FB had no transcription factor binding sequence, including promoter module, matches (GEMS release 3.5.2 June 2003; vertebrate TF binding sequence families (core similarity: 0.75 / matrix similarity: opt); and promoter modules (default parameters: optimized threshold or 80% of maximum score)), while hhyg-3FB had 4 transcription factor binding sequence matches remaining but no promoter modules (Table 10). The following transcription factor binding sequences were found in hhyg-3FB:

1) V\$MINI

Family: Muscle Initiators (2 members)

20 Best match: Muscle Initiator Sequence 1

Ref: Laura L. Lopez & James W. Fickett "Muscle-Specific Regulation of Transcription: A Catalog of Regulatory Elements"

<http://www.cbil.upenn.edu/MTIR/HomePage.html>

25 Position in ORF: -7 to 11

2) V\$PAX5

Family: PAX-5/PAX-9 B-cell-specific activating proteins (4 members)

Best match: B-cell-specific activating protein

Ref: MEDLINE 94010299

30 Position in ORF: 271 to 299

3) V\$AREB

Family: Atp1a1 regulatory element binding (4 members)

Best match: AREB6

Ref: MEDLINE 96061934

Position in ORF: 310 to 322

4) V\$VMYB

Family: AMV-viral myb oncogene (2 members)

5 Best match: v-Myb

Ref: MEDLINE 94147510

Position in ORF: 619 to 629

- Other sequences remaining in hneo-3F included one *E. coli* RBS 8 bases upstream of Met (ORF position 334 to 337); hneo-3FB included a splice acceptor site (+) and *Pst*I site as part of a 5' cloning site for *Sbf*I, and one *E. coli* RBS 8 bases upstream of Met (ORF position 334 to 337); hhyg-3F had no other sequence matches; and hhyg-3FB included a splice acceptor site (+) and *Pst*I site as part of a 5' cloning site for *Sbf*I.
- Subsequently, regulatory sequences were removed from hneo-3F and hhyg-3F yielding hneo-4 and hhyg-4. Then regulatory sequences were removed from hneo-4 yielding hneo-5.

Table 14

| Gene name | TF binding sequences | Promoter modules |
|-----------|----------------------|-------------------|
| | 5' F / ORF / 3' F | 5' F / ORF / 3' F |
| Neo | -- / 53 / -- | -- / 0 / -- |
| hneo-F | 1 / 61 / 2 | 0 / 2 / 0 |
| hneo-3F | 0 / 0 / 0 | 0 / 0 / 0 |
| hneo-3FB | 0 / 0 / 0 | 0 / 0 / 0 |
| Hyg | -- / 74 / -- | -- / 3 / -- |
| hhyg-F | 1 / 94 / 1 | 0 / 4 / 0 |
| hhyg-3F | 1 / 3 / 0 | 0 / 0 / 0 |
| hhyg-3FB | 1 / 3 / 0 | 0 / 0 / 0 |

*Promoter modules are defined as a composite regulatory element, with 2 transcription factor binding sites separated by a spacer, which has been shown to exhibit synergistic or antagonistic function.

Table 15 summarizes the identity of various genes.

Table 15

Pairwise identity of different gene versions

5

Comparisons were of open reading frames (ORFs).

| | neo | hneo | hneo-3 | hneo-4 | hneo-5 | Final hNeo |
|------------|-----|------|--------|--------|--------|------------|
| Neo | -- | 79 | 78 | 78 | 78 | 77 |
| hneo | | -- | 90 | 90 | 90 | 89 |
| hneo-3 | | | -- | 100 | 99 | 98 |
| hneo-4 | | | | -- | 99 | 98 |
| hneo-5 | | | | | -- | 99 |
| Final hNeo | | | | | | -- |

| | hyg | hhyg | hhyg-3 | hHygro | hhyg-4 | Final hHyg |
|------------|-----|------|--------|--------|--------|------------|
| Hyg | -- | 79 | 78 | 73 | 76 | 78 |
| hhyg | | -- | 88 | 83 | 86 | 88 |
| hhyg-3 | | | -- | 94 | 96 | 98 |
| hHygro | | | | -- | 96 | 94 |
| hhyg-4 | | | | | -- | 97 |
| Final hHyg | | | | | | -- |

| Percent Identity | | | | | |
|------------------|---|------|------|---|-----------------------------|
| Divergence | | 1 | 2 | | |
| | 1 | | 82.2 | 1 | Synthetic puro-SEQ ID NO:11 |
| | 2 | 19.6 | | 2 | Starting puro-SEQ ID NO:15 |
| | | 1 | 2 | | |

10

An expression cassette (hNeo-cassette) with a synthetic neomycin gene flanked by a SV40 promoter and a synthetic poly(A) site is shown below.

GGATCCGTTTGC GTATTGGGCGCTCTTCCGCTGATCTGCGCAGCACCA
TGGCCTGAAATAACCTCTGAAAGAGGAACTTGGTTAGCTACCTTCTG

AGGCGGAAAGAAACAGCTGTGGAATGTGTGTCAGTTAGGGTGTGGAA
 AGTCCCCAGGCTCCCCAGCAGGCAGAAAGTATGCAAAGCATGCATCTC
 AATTAGTCAGCAACCAGGTGTGGAAAGTCCCCAGGCTCCCCAGCAGG
 CAGAAGTATGCAAAGCATGCATCTCAATTAGTCAGCAACCATAGTCC
 5 CGCCCCTAACTCCGCCCATCCCGCCCCTAACTCCGCCCAGTTCCGCCC
 ATTCTCCGCCCCATGGCTGACTAATTTTTTTTATTTATGCAGAGGCCG
 AGGCCGCCTCTGCCTCTGAGCTATTCCAGAAGTAGTGAGGAGGCTTT
 TTTGGAGGCCTAGGCTTTTGCAAAAAGCTCGATTCTTCTGACACTAGC
 GCCACCATGATCGAACAAGACGGCCTCCATGCTGGCAGTCCCGCAGC
 10 TTGGGTGGAACGCTTGTTTCGGGTACGACTGGGCCCAGCAGACCATCG
 GATGTAGCGATGCGGCCGTGTTCCGTCTAAGCGCTCAAGGCCGGCCC
 GTGCTGTTTCGTGAAGACCGACCTGAGCGGCCCTGAACGAGCTTCA
 AGACGAGGCTGCCGCCTGAGCTGGCTGGCCACCACCGGCGTACCCT
 GCGCCGCTGTGTTGGATGTTGTGACCGAAGCCGGCCGGGACTGGCTG
 15 CTGCTGGGCGAGGTCCCTGGCCAGGATCTGCTGAGCAGCCACCTTGC
 CCCCCTGAGAAAGTTTCTATCATGGCCGATGCAATGCGGCGCCTGC
 ACACCCTGGACCCCGCTACCTGCCCCCTTCGACCACCAGGCTAAGCAT
 CGGATCGAGCGTGCTCGGACCCGCATGGAGGCCGGCCTGGTGGACCA
 GGACGACCTGGAAGAGGAGCATCAGGGCCTGGCCCCCGCTGAACTGT
 20 TCGCCCGACTGAAAGCCCGCATGCCGGACGGTGAGGACCTGGTTGTC
 ACACACGGAGATGCCTGCCTCCCTAACATCATGGTCGAGAATGGCCG
 CTTCTCCGGCTTCATCGACTGCGGTGCGCTAGGAGTTGCCGACCGCTA
 CCAGGACATCGCCCTGGCCACCCGCGACATCGCTGAGGAGCTTGGCG
 GCGAGTGGGCCGACCGCTTCTTAGTCTTGTACGGCATCGCAGCTCCC
 25 GACAGCCAGCGCATCGCCTTCTACCGCTTGCTCGACGAGTTCTTTTAA
 TGATCTAGAACCGGTCATGGCCGCAATAAAATATCTTTATTTTCATTA
 CATCTGTGTGTTGGTTTTTTGTGTGTTTCGAACTAGATGCTGTGCGAC
 (SEQ ID NO:44).

30 An expression cassette (hPuro-cassette) with a synthetic puromycin gene flanked
 by a SV40 promoter and a synthetic poly(A) site is shown below.

GGATCCGTTTGC GTATTGGGCGCTCTTCCGCTGATCTGCGCAGCACCA
 TGGCCTGAAATAACCTCTGAAAGAGGAACTTGGTTAGCTACCTTCTG

AGGCGGAAAGAACCAGCTGTGGAATGTGTGTCAGTTAGGGTGTGGAA
 AGTCCCCAGGCTCCCCAGCAGGCAGAAGTATGCAAAGCATGCATCTC
 AATTAGTCAGCAACCAGGTGTGGAAAGTCCCCAGGCTCCCCAGCAGG
 CAGAAGTATGCAAAGCATGCATCTCAATTAGTCAGCAA.CCATAGTCC
 5 CGCCCCTAACTCCGCCCATCCCGCCCCTAACTCCGCCCAGTTCCGCCC
 ATTCTCCGCCCCATGGCTGACTAATTTTTTTTATTTATGCAGAGGCCG
 AGGCCGCCTCTGCCTCTGAGCTATTCCAGAAGTAGTGA.GGAGGCTTT
 TTTGGAGGCCTAGGCTTTTGCAAAAAGCTCGATTCTTCTGACACTAGC
 GCCACCATGACCGAGTACAAGCCTACCGTGCGCCTGGCCACTCGCGA
 10 TGATGTGCCCCGCGCCGTCCGCACTCTGGCCGCCGCTTTCGCCGACTA
 CCCCCTACCCGGCACACCGTGACCCCGACCGGCACATCGAGCGTG
 TGACAGAGTTGCAGGAGCTGTTCTGACCCGCGTCGGGCTGGACATC
 GGCAAGGTGTGGGTAGCCGACGACGGCGCGGCCGTGG.CCGTGTGGA
 CTACCCCCGAGAGCGTTGAGGCCGGCGCCGTGTTCCGCCGAGATCGGC
 15 CCCCGAATGGCCGAGCTGAGCGGCAGCCGCCTGGCCGCCAGCAGCA
 AATGGAGGGCCTGCTTGCCCCCATCGTCCCAAGGAGCCTGCCTGGT
 TTCTGGCCACTGTAGGAGTGAGCCCCGACCACCAGGGCAAGGGCTTG
 GGCAGCGCCGTCTGTGTTGCCCGGCGTAGAGGCCGCCGAACGCGCCGG
 TGTGCCCGCCTTTCTCGAAACAAGCGCACCAAGAAACCTTCCATTCTA
 20 CGAGCGCCTGGGCTTCACCGTGACCGCCGATGTCGAGGTGCCCGAGG
 GACCTAGGACCTGGTGTATGACACGAAAACCTGGCGCCTAATGATCT
 AGAACCGGTCATGGCCGCAATAAAATATCTTTATTTTC.ATTACATCTG
 TGTGTTGGTTTTTTGTGTGTTTGAAGTACTAGATGCTGTGAC (SEQ ID
 NO:11);

25

hpuro:

GCTAGCGCCACCATGACCGAGTACAAGCCCACCGTGCGCCTGGCCAC
 CCGCGACGACGTGCCCCGCGCCGTGCGCACCCCTGGCCGCCGCCTTCG
 CCGACTACCCCGCCACCCGCCACACCGTGACCCCGACCGCCACATC
 30 GAGCGCGTGACCGAGCTGCAGGAGCTGTTCTGACCCGCGTGCGCCT
 GGACATCGGCAAGGTGTGGGTGGCCGACGACGGCGCCGCCGTGGCC
 GTGTGGACCACCCCGAGAGCGTGAGGCCGGCGCCGTGTTCCGCCGA
 GATCGGCCCCCGCATGGCCGAGCTGAGCGGCAGCCGCCTGGCCGCC

AGCAGCAGATGGAGGGCCTGCTGGCCCCCACC GCCCCAAGGAGCCC
 GCCTGGTTCTCTGGCCACCGTGGGCGTGAGCCCCGACCACCAGGGCAA
 GGGCCTGGGCAGCGCCGTGGTGCTGCCCCGGCGTGGAGGCCGCCGAGC
 GCGCCGGCGTGCCCGCCTTCTGGAGACCAGCGCCCCCGCAACCTG
 5 CCCTTCTACGAGCGCCTGGGCTTCACCGTGACCGCCGACGTGGAGGT
 GCCCGAGGGCCCCCGCACCTGGTGCATGACCCGCAAGCCCGGCGCCT
 AATGATCTAGA (SEQ ID NO:91);

hpuro-1:

10 gctagcgccaccatgaccgagtacaagcctaccgtgcgccctggccactcgcgatgatgtgccccgcgccgtccgc
 actctggccgcccgtttcgccgactaccccgctacccggcacaccgtggaccccgaccggcacatcgagcgtgtg
 acagagttgcaggagctgttctgacccgctcgggctggacatcggaagggtgtgggtagccgacgacggcgc
 ggccgtggccgtgtggactaccccgagagcgttgaggccggcgccgtgttcgccgagatcgccccgaatgg
 ccgagctgagcggcagccgcctggccgccagcagcaaatggaggcctgcttgcccccatcgtcccaaggag
 15 cccgcctggtttctggccactgtaggagttagccccgaccaccagggaaggccttgggcagcgccgtcgtgtg
 cccggcgtagaggccgccgaacgcgggtgtgcccgcctttctggagacaagcgctccgcgtaacctccattct
 acgagcgccctggccttaccgtgaccgcccgatgtcgaggtgcccgagggaacccggacctggtgcatgactcgc
 aagcctggcgccctaatgatctaga (SEQ ID NO:92); and

20 hpuro-2

GCTAGCGCCACCATGACCGAGTACAAGCCTACCGTGCGCCTGGCCAC
 TCGCGATGATGTGCCCCGCGCCGTCCGCACTCTGGCCGCCGCTTTCGC
 CGACTACCCCGCTACCCGGCACACCGTGGACCCCGACCGGCACATCG
 AGCGTGTGACAGAGTTGCAGGAGCTGTTCTTGACCCGCGTCGGGCTG
 25 GACATCGGCAAGGTGTGGGTAGCCGACGACGGCGCGGCCGTGGCCG
 TGTGGACTACCCCGAGAGCGTTGAGGCCGGCGCCGTGTTCCGCCGAG
 ATCGGCCCCCGAATGGCCGAGCTGAGCGGCAGCCGCCTGGCCGCCCA
 GCAGCAAATGGAGGGCCTGCTTGCCCCCATCGTCCCAAGGAGCCTG
 CCTGGTTTCTGGCCACTGTAGGAGTGAGCCCCGACCACCAGGGCAAG
 30 GGCTTGGGCAGCGCCGTCGTGTTGCCCGGCGTAGAGGCCGCCGAACG
 CGCCGGTGTGCCCCGCTTTCTCGAAACAAGCGCACCAAGAAACCTTC
 CATTCTACGAGCGCCTGGGCTTCACCGTGACCGCCGATGTTCGAGGTG
 CCCGAGGGACCTAGGACCTGGTGTATGACACGAAAACCTGGCGCCTA

ATGATCTAGA (SEQ ID NO:93).

The starting puro sequence (from psi STRIKE) has SEQ ID NO:15

(atgaccgagt acaagcccac ggtgcgcctc gccacccgcg acgacgtccc ccggggccgta
 5 cgcaccctcg ccgccgcgtt cgccgactac cccgccacgc gccacaccgt cgacccggag
 cgccacatcg agcgggtcac cgagctgcaa gaactcttcc tcacgcgcgt cgggctcgac
 atcggaagg tgtgggtcgc ggacgacggc gccgcggtgg cggcttgac cacgccggag
 agcgtgaag cgggggcggt gtccgccgag atcgcccgcc gcatggccga gttgagcgt
 tcccggtg ccgcgcagca acagatggaa ggctctctgg cccgcaccg gccaaggag
 10 cccgcgtggt tcctggccac cgtcggcgtg tcgccgacc accagggcaa gggctgggc
 agcgccctcg tgctccccg agtgaggcg gccgagcgcg ccggggtgcc cgccttctg
 gagacctcg cccccgcaa cctcccctc tacgagcggc tcggcttac cgtaccgcc
 gacgtcagg tgcccgaagg accgcgcacc tggatcatga cccgaagcc cgggtgcc).

15 Other synthetic *hyg* and *neo* genes include

hneo-1:

CCACTCAGTGGCCACCATGATCGAGCAGGACGGCCTcCAtGcTGGCAGt
 CCCGCaGCCTGGGTcGAGCGCtTGTTcGGgTACGACTGGGCCCAGCAG
 ACCATCGGaTGtAGCGAtGCCGCaGTGTTCCGCCTGAGCGCtCAaGGCCG
 20 gCCCGTGCTGTTCGTGAAGACCGACCTGAGCGGCGCCCTGAACGAGC
 TtCAaGACGAGGcTCCCCGCCTGAGCTGGCTGGCCACCACCGGtGTaCC
 CTGCGCCGcTGTgTGGAgtGTgTGACCGAaGCCGGCCGCGACTGGCTGC
 TGCTGGGCGAGGTGCCtGGCCAGGACCTGCTGAGCAGCCACCTGGCC
 CCCGcTGAaAGGTGAGCATCATGGCCGACGCCATGCGgCGCCTGCAC
 25 ACCCTGGACCCCGcTACaTGCCCCCTTCGACCACCAGGcTAAGCACCGC
 ATCGAGCGgGCtCGgACCCGCATGGAGGCCGGCCTGGTGGACCAGGAC
 GACCTGGACGAGGAGCACCAGGGCCTGGCCCCCGcTGAaCTGTTCCGC
 CGCCTGAaAGCCGCATGCCgGACGGtGAGGACCTGGTtGTGACaCACG
 GCGACGCCTGCCTcCCtAACATCATGGTcGAGAACGGgCGCTTcTcCGG
 30 TTCATCGACTGCGGCCGCCTGGGCGTtGCCGACCGCTACCAGGACATC
 GCCCTGGCCACCCGCGACATCGCCGAGGAGCTGGGCGGCGAGTGGG
 CCGACCGCTTCCTGGTctTGtACGGCATCGCaGCtCCCGACAGCCAGCG
 CATCGCCTTCTACCGCCTGCTGGACGAGTTCTTCTAgTAACCAGgCTCT

GG (SEQ ID NO:38);

hneo-2

CCACTCcGTGGCCACCATGATCGAaCAaGACGGCCTcCAtGcTGGCAGtC
5 CCGCaGcTGGGTcGAaCGcTGTTCGGgTACGACTGGGCCCAGCAGAC
CATCGGaTgtAGCGAtGCgGCCGTGTTCCGtCTaAGCGCtCAaGGCCGgCC
CGTGCTGTTCTGTGAAGACCGACCTGAGCGGCGCCCTGAACGAGCTtCA
aGACGAGGcTCCCCGCCTGAGCTGGCTGGCCACCACCGGtGTaCCCTGC
GCCGcTGTgtTGGAtGTtGTGACCGAaGCCGGCCGgGACTGGCTGCTGCT
10 GGGCGAGGTcCCtGGCCAGGAiCTGCTGAGCAGCCACCTtGCCCCGcT
GAGAAGGTttcCATCATGGCCGAtGCaATGCGgCGCCTGCACACCCTGG
ACCCCGcTACaTGCCCCCTTCGACCACCAGGcTAAGCAtCGgATCGAGCGt
GcTCGgACCCGCATGGAGGCCGGCCTGGTGGACCAGGACGACCTGGA
CGAGGAGCAtCAGGGCCTGGCCCCCGcTGAaCTGTTCGCCCCGCTGAaA
15 GCCCCGCATGCCgGACGGtGAGGACCTGGTtGTGACaCAtGGaGAtGCCTG
CCTcCCtAACATCATGGTcGAGAAtGGcCGCTTcTcCGGCTTCATCGACTG
CGGtCGCCTaGGaGTtGCCGACCGCTACCAGGACATCGCCCTGGCCACC
CGCGACATCGcTgAGGAGCTtGGCGGCGAGTGGGCCGACCGCTTcTtTaG
TctTGTACGGCATCGCaGcTCCCGACAGCCAGCGCATCGCCTTCTACCG
20 CCTGCTcGACGAGTTCTTtTAATGACCAGgCTCTGG (SEQ ID NO:39);

hhyg-1

CCACTCAGTGGCCACCATGAAGAAGCCCGAGCTGACCGCTACCAGCG
TTGAGAAGTTCCTGATCGAGAAGTTCGACAGCGTGAGCGACCTGATG
CAGTTAAGCGAGGGCGAGGAAAGCCGCGCCTTCAGCTTCGATGTGCG
25 CGGACGCGGCTATGTACTGCGGGTGAATAGCTGCGCTGATGGCTTCT
ACAAAGACCGCTACGTGTACCGCCACTTCGCCAGCGCTGCACTGCCC
ATCCCCGAGGTGCTGGACATCGGCGAGTTCAGCGAGAGCCTGACATA
CTGCATCAGCCGCCGCGCTCAAGGCGTGACTCTCCAAGACCTGCCCCG
AGACAGAGCTGCCCCGCTGTGCTACAGCCTGTCGCCGAGGCTATGGAC
30 GCTATTGCCGCCGCCGACCTGAGCCAGACCAGCGGCTTCGGCCCATT
CGGGCCCCAAGGCATCGGCCAGTACACCACCTGGCGCGACTTCATCT
GCGCCATTGCTGATCCCCATGTCTACCACTGGCAGACCGTGATGGAC
GACACCGTGAGCGCCAGCGTAGCTCAAGCCCTGGACGAGCTGATGCT

GTGGGCCGAGGACTGCCCCGAGGTGCGCCATCTCGTCCATGCCGACT
TCGGCAGCAACAACGTCCTGACCGACAACGGCCGCATCACCGCCGTA
ATCGACTGGAGCGAGGCCATGTTCTGGGGACAGTCAGTACGAGGTGGC
CAACATCTTCTTCTGGCGGCCCTGGCTGGCCTGCATGGAGCAGCAAA
5 CCCGCTACTTCGAGCGCCGCCATCCCGAGCTGGCCGGCAGCCCCCGT
CTGCGAGCCTACATGCTGCGCATCGGCCTGGATCAGCTCTACCAGAG
CCTCGTGGACGGCAACTTCGACGATGCTGCCTGGGCTCAAGGCCGCT
GCGATGCCATCGTCCGCAGCGGGGCCGGCACCGTCGGTCGCACACAA
ATCGCTCGCCGGAGCGCCGCCGTATGGACCGACGGCTGCGTCGAGGT
10 GCTGGCCGACAGCGGCAACCGCCGGCCCAGTACACGACCGCGCGCTA
AGGAGTAGTAACCAGCTCTTGG (SEQ ID NO:30);

hhyg-2:

CCACTCCGTGGCCACCATGAAAGAAGCCCGAGCTGACCGCTACCAGCG
15 TTGAAAAATTTCTCATCGAGAGTTCGACAGTGTGAGCGACCTGATG
CAGTTGTCGGAGGGCGAAGAGAGCCGAGCCTTCAGCTTCGATGTCGG
CGGACGCGGCTATGTACTGCGGGTGAATAGCTGCGCTGATGGCTTCT
ACAAAGACCGCTACGTGTACCGCCACTTCGCCAGCGCTGCACTACCC
ATCCCCGAAGTGTGGACATCGGCGAGTTCAGCGAGAGCCTGACATA
20 CTGCATCAGTAGACGCGCCCAAGGCGTTACTCTCCAAGACCTCCCCG
AAACAGAGCTGCCTGCTGTGTTACAGCCTGTCGCCGAAGCTATGGAT
GCTATTGCCGCCGCCGACCTCAGTCAAACCAGCGGCTTCGGCCCATTT
CGGGCCCCAAGGCATCGGCCAGTACACAACCTGGCGGGATTTCATTT
GCGCCATTGCTGATCCCCATGTCTACCACTGGCAGACCGTGATGGAC
25 GACACCGTGTCCGCCAGCGTAGCTCAAGCCCTGGACGAACTGATGCT
GTGGGCCGAAGACTGTCCCGAGGTGCGCCACCTCGTCCATGCCGACT
TCGGCAGCAACAACGTCCTGACCGACAACGGCCGCATCACCGCCGTA
ATCGACTGGAGCGAGGCTATGTTCTGGGGACAGTCAGTACGAGGTGGC
CAACATCTTCTTCTGGCGGCCCTGGCTGGCTTGCATGGAGCAGCAGA
30 CTCGCTACTTCGAGCGCCGGCATCCCGAGCTGGCCGGCAGCCCTCGT
CTGCGAGCCTACATGCTGCGCATCGGCCTGGATCAGCTCTACCAGAG
CCTCGTGGACGGCAACTTCGACGATGCTGCCTGGGCTCAAGGCCGCT
GCGATGCCATCGTCCGCAGCGGGGCCGGCACCGTCGGTCGCACACAA

ATCGCTCGCCGGAGCGCCGCGGTATGGACCGACGGCTGCGTCGAGGT
 GCTGGCCGACAGCGGCAACCGCCGGCCAGTACACGACCGCGCGCTA
 AGGAGTAGTAACCAGCTCTTGG (SEQ ID NO:42);

- 5 hHygro (*Sac*I site in ORF near 5' end, insert in-frame linker coding for 12 amino acids at 3' end, and *Sna*BI site added at 3' end in ORF)
- aagcttgctagcgccaccatgaagaagcccagctcaccgctaccagcgttgaaaaatttctcatcgagaagttcga
 cagtgtgagcgacctgatgcagttgtcggaggcggaagagagccgagccttcagcttcgatgtcggcgagcgcg
 ctatgtactgcggtgaatagctgcgctgatggcttctacaaagaccgctacgtgtaccgccacttcgccagcgctgc
 10 actaccatccccgaagtgttgacatcggcgagttcagcgagacgtgacatactgcatcagtagacgcgccc
 ggcgttactctccaagacctccccgaaacagagctgcctgtgttacagcctgtcgcggaagctatggatgctatt
 gccgccggcgacctcagtcacaccagcggttcggccattcgggcccccaaggcatcgccagtagacacaacctg
 gcgggatttcatttgcgcattgctgatccccatgtctaccactggcagaccgtgatggacgacaccgtgtccgccag
 cgtagctcaagccctggacgaactgatgtgtggccgaagactgtcccaggtgcgccacctgtccatgccgac
 15 ttggcgagcaacaacgtcctgaccgacaacggccgcatcaccgccgaatcgactgtccgaagctatgttcgggg
 acagttagtagaggtggccaacatcttcttggcgccctggctggcttgcatggagcagcagactcgtacttc
 gagcgccggcatcccgagctggccggcagccctcgtctgcgagcctacatgtgcgcatcggcctggatcagctc
 taccagagcctcgtggacggcaacttcgacgatgtgcctgggctcaaggccgctgcgatgccatcgtccgcagc
 ggggcccggcaccgtcgtcgacacaaaatcgtcggcgagcgccgctatggaccgacggctgcgtcgaggtg
 20 gctggccgacagcggaaccggcccgccagtagacgaccgcgcgctaaggagggtggcgaggaggagcggtgg
 cggaggttctacgtatagctagactcgag (SEQ ID NO:70);

hhyg-4

- atgaagaagcccagctcaccgctaccagcgttgaaaaatttctcatcgagaagttcgacagtgtgagcgacctgat
 25 gcagttgtcggaggcggaagagagccgagccttcagcttcgatgtcggcgagcgcggctatgtactgcggtgaa
 tagctgcgctgatggcttctacaaagaccgctacgtgtaccgccacttcgccagcgcgcactaccatccccgaag
 tgttgacatcggcgagttcagcgagagcctgacatactgcatcagtagacgcgcccaaggcggttactctccaaga
 cctccccgaaacagagctgcctgtgtgttacagcctgtcggcgaagctatggatgctattgccgccggcgacctca
 gtcaaaccagcggttcggccattcgggcccccaaggcatcgccagtagacacaacctggcgggatttcatttgcgc
 30 cattgtgatccccatgtctaccactggcagaccgtgatggacgacaccgtgtccgcagcgtagctcaagccctgg
 acgaactgatgtgtggccgaagactgtcccaggtgcgccacctgtccatgccgacttcggcagcaacaacgt
 cctgaccgacaacggccgcatcaccgccgaatcgactggtccgaagctatgttcgggacagtcagtagaggtg
 gccaacatcttcttggcgccctggctggcttgcatggagcagcagactcgtacttcgagcgccggcatccccga

gctggccggcagccctcgtctgcgagcctacatgctgcgcacatggcctggatcagctctaccagagcctcgtggac
ggcaacttcgacgatgctgcctgggctcaaggccgctgcgatgccatcgtccgcagcggggccggcaccgtcgg
cgcacacaaatcgctcgccggagcgcagccgtagtgaccgacggctgcgtcgagggtgctggccgacagcggca
accgccggcccagtagcacgaccgcgcgctaaggaaggcggaggtagtggtggcggaggtagctacgta

5 (SEQ ID NO:71);

hneo-4:

GCTAGCGCCACCATGATCGAACAAGACGGCCTCCATGCTGGCAGTCC
CGCAGCTTGGGTTCGAACGCTTGTTTCGGGTACGACTGGGCCCAGCAGA
10 CCATCGGATGTAGCGATGCGGCCGTGTTCCGTCTAAGCGCTCAAGGC
CGGCCCCGTGCTGTTTCGTGAAGACCGACCTGAGCGGCGCCCTGAACGA
GCTTCAAGACGAGGCTGCCCCGCTGAGCTGGCTGGCCACCACCGGTG
TACCCTGCGCCGCTGTGTTGGATGTTGTGACCGAAGCCGGCCGGGAC
TGGCTGCTGCTGGGCGAGGTCCCTGGCCAGGATCTGCTGAGCAGCCA
15 CCTTGCCCCCGCTGAGAAGGTTTCCATCATGGCCGATGCAATGCGGC
GCCTGCACACCCTGGACCCCGCTACATGCCCCTTCGACCACCAGGCT
AAGCATCGGATCGAGCGTGCTCGGACCCGCATGGAGGCCGGCCTGGT
GGACCAGGACGACCTGGACGAGGAGCATCAGGGCCTGGCCCCCGCT
GAACTGTTTCGCCCCGCTGAAAGCCCGCATGCCGGACGGTGAGGACCT
20 GGTTGTGACACATGGTGATGCCTGCCTCCCTAACATCATGGTTCGAGA
ATGGCCGCTTCTCCGGCTTCATCGACTGCGGTTCGCTAGGAGTTGCCG
ACCGCTACCAGGACATCGCCCTGGCCACCCGCGACATCGCTGAGGAG
CTTGGCGGCGAGTGGGCGACCGCTTCTTAGTCTTGTACGGCATCGC
AGCTCCCGACAGCCAGCGCATCGCCTTCTACCGCCTGCTCGACGAGT
25 TCTTTTAATCTAGA

(SEQ ID NO:72);

and

hneo-5:

GCTAGCGCCACCATGATCGAACAAGACGGCCTCCATGCTGGCAGTCC
30 CGCAGCTTGGGTTCGAACGCTTGTTTCGGGTACGACTGGGCCCAGCAGA
CCATCGGATGTAGCGATGCGGCCGTGTTCCGTCTAAGCGCTCAAGGC
CGGCCCCGTGCTGTTTCGTGAAGACCGACCTGAGCGGCGCCCTGAACGA
GCTTCAAGACGAGGCTGCCCCGCTGAGCTGGCTGGCCACCACCGGCG

TACCCTGCGCCGCTGTGTTGGATGTTGTGACCGAAGCCGGCCGGGAC
 TGGCTGCTGCTGGGCGAGGTCCCTGGCCAGGATCTGCTGAGCAGCCA
 CCTTGCCCCCGCTGAGAAGGTTTCTATCATGGCCGATGCAATGCGGC
 GCCTGCACACCCTGGACCCCGCTACCTGCCCCTTCGACCACCAGGCT
 5 AAGCATCGGATCGAGCGTGCTCGGACCCGCATGGAGGCCGGCCTGGT
 GGACCAGGACGACCTGGACGAGGAGCATCAGGGCCTGGCCCCCGCT
 GAACTGTTGCCCCGACTGAAAGCCCGCATGCCGGACGGTGAGGACCT
 GGTTGTCACACACGGAGATGCCTGCCTCCCTAACATCATGGTCGAGA
 ATGGCCGCTTCTCCGGCTTCATCGACTGCGGTCGCCTAGGAGTTGCCG
 10 ACCGCTACCAGGACATCGCCCTGGCCACCCGCGACATCGCTGAGGAG
 CTTGGCGGCGAGTGGGCCGACCGCTTCTTAGTCTTGTACGGCATCGC
 AGCTCCCGACAGCCAGCGCATCGCCTTCTACCGCTTGCTCGACGAGTT
 CTTTAAATGATCTAGA(SEQ ID NO:73).

15 The synthetic nucleotide sequence of the invention may be employed in
 fusion constructs. For instance, a synthetic sequence for a selectable polypeptide
 may be fused to a wild-type sequence or to another synthetic sequence which
 encodes a different polypeptide. For instance, the *neo* sequence in the following
 examples of a synthetic *Renilla* luciferase-*neo* sequence may be replaced with a
 20 synthetic *neo* sequence of the invention:
 atggcttcaaggtgtacgaccccgagcaacgcaaacgcatgatcactgggcctcagtgggtgggctcgtgcaagc
 aaatgaacgtgctggactcctcatcaactactatgattccgagaagcacgcccagagaacgccgtgattttctcatggt
 taacgctgcctccagctacctgtggaggcacgtcgtgcctcacatcgagcccgtggctagatgcatcatccctgatct
 gatcgggaatgggtaagtccggcaagagcgggaatggctcatatcgccctcctggatcactacaagtacctaccgctt
 25 ggttcgagctgctgaaccttccaaagaaaatcatctttgtggccacgactgggggcttgtctggcctttcactactc
 ctacgagcaccaagacaagatcaaggccatcgtccatgctgagagtgctgtggacgtgatcgagtcctgggacga
 gtggcctgacatcgaggaggatcgcctgatcaagagcgaagagggcgagaaaatggtgcttgagaataacttc
 ttctcgagaccatgctccaagcaagatcatcggaactggagcctgaggagttcgtgcctacctggagccatt
 caaggagaagggcgaggttagacggcctaccctctcctggcctcgcgagatccctctcgttaagggaggcaagcc
 30 cgacgtcgtccagattgtccgcaactacaacgctacctcgggcccagcgacgatctgcctaagatgttcacgagtc
 cgaccctgggttctttccaacgctattgtcgaggagctaagaagttccctaaccaggagttcgtgaaggtgaaggg
 cctccacttcagccaggaggagcgtccagatgaaatgggtaagtacatcaagagcttcgtggagcgcgtgctgaag
 aacgagcagaccgggtggtgggagcggaggtggcgatcaggtggcgaggctccggagggtgaacaagatg

gattgcacgcaggttctccggccgcttgggtggagaggctattcggctatgactgggcacaacagacaatcggtg
ctctgatgccgccgtgttccggctgtcagcgcagggcgcccggttcttttgaagaccgacctgtccgggtgccct
gaatgaactgcaggacgaggcagcgcggctatcgtggctggccacgacggcggttccttgcgcagctgtgctcga
cgttgtcactgaagcgggaagggactggctgctattggcggaagtgccggggcaggatctcctgtcatctcaccttg
5 ctcctgccgagaaagtatccatc atggctgatgcaatcgggcggtgcatacgttgatccggctacctgccattcg
accaccaagcgaacatcgcatcgagcgcagcacgtactcggatggaagccggcttctgtcatcaggatgatctgga
cgaagagcatcaggggctcgcgccagccgaactgttcgccaggctcaaggcgcgatccccgacggcgaggat
ctcgtcgtgacctatggcgatgcctgttggcgaatatcatggtgaaaatggccgcttttctgattcatcgactgtg
gccggctgggtgtggcggaccgctatcaggacatagcgttggctacctgtatattgctgaagagcttggcggcga
10 atgggctgaccgcttctcgtgctttacggctatcgccgctcccgattcgcagcgcacatcgcccttatcgcccttctgacg
agttcttctaa (hrl-neo fusion; SEQ ID NO:12)
and
atgattgaacaagatggattgcacgcaggttctccggccgcttgggtggagaggctattcggctatgactgggcaca
acagacaatcggtgctctgatgccgccgtgtccggctgtcagcgcagggcgcccggttcttttgaagaccg
15 acctgtccgggtgccctgaatgaa.ctgcaggacgaggcagcgcggctatcgtggctggccacgacggcggttcctt
gcgcagctgtgctcagctgtc actgaagcgggaagggactggctgctattggcggaagtccggggcaggatc
tcctgtcatctcaccttgcctcgcgcgagaaagtatccatcatggctgatgcaatcgggcggtgcatacgttgatcc
ggctacctgccattcgaccacc aagcgaacatcgcatcgagcgcagcgtactcggatggaagccggcttctgtc
gatcaggatgatctggacgaagagcatcaggggctcgcgccagccgaactgttcgccaggctcaaggcgcgat
20 gcccgcggcgaggatctgtcgtgacctatggcgatgcctgttggcgaatatcatggtgaaaatggccgctttt
ctggattcatcgactgtggccggctgggtgtggcggaccgctatcaggacatagcgttggctacctgtatattgctg
aagagcttggcggcgaatggcctgaccgcttctcgtgctttacggatcgcgcgctcccgattcgcagcgcacatgcc
ttctatcgcccttctgacgagttcttcaccgggtgtgggagcggaggtggcggatcagggtggcggaggctccggag
gggcttcaagggtgtacgacctc gagcaacgcaaacgcgatgactgggcctcagtggtgggctcgtgcaagc
25 aatgaacgtgctggactccttcataactactatgattccgagaagcacgccgagaacgccgtgattttctcatggt
taacgctgcctccagctacctgtggaggcacgtcgtgcctcacatcgagcccggtggctagatgcatcatccctgatct
gatcggaatgggtaagtccggcaagagcgggaatggctcatatcgccctcctggatcactacaagtacctcaccgctt
ggttcgagctgctgaaccttcaaaagaaaatcatcttttggggccacgactggggggcttctgtgcctttcactactc
ctacgagcaccaagacaagatcaaggccatcgtccatgctgagagtgctgtggacgtgatcagtcctgggacga
30 gtggcctgacatcgaggaggataatcgccctgatcaagagcgaagaggcgagaaaatgggtgcttgagaataacttc
ttcgtcgagacctgctccaagcaagatcatcgggaaactggagcctgaggagttcgtgcctacctggagccatt
caaggagaagggcgagggttagcggcctaccctctcctggcctcgcgagatccctctctgtaaggagggaagcc
cgacgtcgtccagattgtccgcaactacaacgcctaccttcgggccagcgacgatctgcctaagatgttcatcgagtc

cgaccctgggttctttccaacgctattgtcgagggagctaagaagttccctaacaccgagttcgtgaaggtgaaggg
 cctccattcagccaggaggacgctccagatgaaatgggtaagtacatcaagagcttcgtgagcgcgtgctgaag
 aacgagcagtaa (neo-hrl-fusion; SEQ ID NO:13).

5

Example 5

Transcription Factor Binding Sites Used to Identify Sites
in Selected Synthetic Sequences

TF binding site libraries

The TF binding site library ("Matrix Family Library") is part of the
 10 GEMS Launcher package. Table 16 shows the version of the Matrix Family
 Library which was used in the design of a particular sequence and Table 17
 shows a list of all vertebrate TF binding sites ("matrices") in Matrix Family
 Library Version 2.4, as well as all changes made to vertebrate matrices in later
 versions up to 4.1 (section "GENOMATIX MATRIX FAMILY LIBRARY
 15 INFORMATION Versions 2.4 to 4.1"). (Genomatix has a copyright to all
 Matrix Library Family information).

Table 16

| Synthetic DNA sequence | Genomatix Matrix Family Library |
|-------------------------------|--|
| pGL4B-NN3* | Version 2.4 May 2002 |
| luc2A8 and luc2B10 | Version 3.0 Nov 2002 Version 3.1.1 April 2003 |
| hhyg3 hneo3 | Version 3.1.2 June 2003 |
| hhyg4 | Version 3.3 August 2003 |
| <i>SpeI-NcoI-Ver2</i> ** | Version 4.0 Nov 2003 |
| hneo5 hpuro2 | Version 4.1 Feb 2004 |

20

**NotI-NcoI* fragment in pGL4 including *amp* gene (pGL4B-NN3)

***SpeI-NcoI-Ver2* (replacement for *SpeI-NcoI* fragment in pGL4B-NN3)

Table 17

GENOMATIX MATRIX FAMILY LIBRARY INFORMATION
Versions 2.4 to 4.1

5 A. Matrix Family Library Version 2.4

Matrix Family Library Version 2.4 (May 2002) contains 412 weight matrices in 193 families

(Vertebrates: 275 matrices in 106 families)

Vertebrates

| Family | Family Information | Matrix Name | Information |
|---------------|---|------------------------|---|
| <u>V\$AHR</u> | AHR-arnt heterodimers and AHR-related factors | <u>V\$AHRARNT.01</u> | aryl hydrocarbon receptor / Arnt heterodimers |
| | | <u>V\$AHR.01</u> | aryl hydrocarbon / dioxin receptor |
| | | <u>V\$AHRARNT.02</u> | aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$AP1</u> | AP1 and related factors | <u>V\$AP1.01</u> | AP1 binding site |
| | | <u>V\$AP1.02</u> | activator protein 1 |
| | | <u>V\$AP1.03</u> | activator protein 1 |
| | | <u>V\$AP1FJ.01</u> | activator protein 1 |
| | | <u>V\$NFE2.01</u> | NF-E2 p45 |
| | | <u>V\$VMAF.01</u> | v-Maf |
| | | <u>V\$TCF11MAFG.01</u> | TCF11/MafG heterodimers, binding to subclass of AP1 sites |
| | | <u>V\$BEL1.01</u> | Bel-1 similar region |
| <u>V\$AP2</u> | Activator Protein 2 | <u>V\$AP2.01</u> | activator protein 2 |
| <u>V\$AP4</u> | AP4 and Related proteins | <u>V\$AP4.01</u> | activator protein 4 |
| | | <u>V\$AP4.02</u> | activator protein 4 |
| | | <u>V\$TH1E47.01</u> | Thing1/E47 heterodimer, TH1 bHLH member specific expression in a variety of embryonic tissues |

| Family | Family Information | Matrix Name | Information |
|----------------|--|----------------------------|--|
| | | <u>V\$TAL1ALPHA E47.01</u> | Tal-1alpha/E47 heterodimer |
| | | <u>V\$TAL1BETA E47.01</u> | Tal-1beta/E47 heterodimer |
| | | <u>V\$TAL1BETA ITF2.01</u> | Tal-1beta/ITF-2 heterodimer |
| | | <u>V\$AP4.03</u> | activator protein 4 |
| <u>V\$AREB</u> | Atp1a1 regulatory element binding | <u>V\$AREB6.04</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| | | <u>V\$AREB6.02</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| | | <u>V\$AREB6.03</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| | | <u>V\$AREB6.01</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| <u>V\$ARP1</u> | Apolipoprotein aI and cIII gene Repressor Protein | <u>V\$ARP1.01</u> | apolipoprotein AI regulatory protein 1 |
| <u>V\$BARB</u> | BARbiturate-Inducible El. box from Pro+eukaryot. genes | <u>V\$BARBIE.01</u> | barbiturate-inducible element |
| <u>V\$BCL6</u> | POZ domain zinc finger expressed in B-Cells | <u>V\$BCL6.01</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| | | <u>V\$BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$BRAC</u> | Brachyury gene, mesoderm | <u>V\$TBX5.01</u> | T-Box factor 5 site (TBX5), mutations related to Holt-Oram |

| Family | Family Information | Matrix Name | Information |
|----------------|---|-----------------------|--|
| | developmental factor | | syndrome |
| | | <u>V\$BRACH.01</u> | Brachyury |
| <u>V\$BRNF</u> | Brn POU domain factors | <u>V\$BRN3.01</u> | POU transcription factor Brn-3 |
| | | <u>V\$BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$CABL</u> | C-abl DNA binding sites | <u>V\$CABL.01</u> | Multifunctional c-Abl src type tyrosine kinase |
| <u>V\$CART</u> | Cart-1 (cartilage homeoprotein 1) | <u>V\$XVENT2.01</u> | Xenopus homeodomain factor Xvent-2; early BMP signaling response |
| | | <u>V\$CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |
| <u>V\$CDXF</u> | Vertebrate caudal related homeodomain protein | <u>V\$CDX2.01</u> | Cdx-2 mammalian caudal related intestinal transcr. factor |
| <u>V\$CEBP</u> | Ccaat/Enhancer Binding Protein | <u>V\$CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| | | <u>V\$CEBP.02</u> | C/EBP binding site |
| <u>V\$CHOP</u> | CHOP binding protein | <u>V\$CHOP.01</u> | heterodimers of CHOP and C/EBPalpha |
| <u>V\$CLOX</u> | CLOX and CLOX homology (CDP) factors | <u>V\$CDPCR3HD.01</u> | cut-like homeodomain protein |
| | | <u>V\$CDP.01</u> | cut-like homeodomain protein |
| | | <u>V\$CDP.02</u> | transcriptional repressor CDP |
| | | <u>V\$CDPCR3.01</u> | cut-like homeodomain protein |
| | | <u>V\$CLOX.01</u> | Clox |
| <u>V\$CMYB</u> | C-MYB, cellular transcriptional activator | <u>V\$CMYB.01</u> | c-Myb, important in hematopoiesis, cellular equivalent to avian myoblastosis virus |

| Family | Family Information | Matrix Name | Information |
|----------------|---|-------------------------|--|
| | | | oncogene v-myb |
| <u>V\$COMP</u> | factors which COoperate with Myogenic Proteins | <u>V\$COMP1.01</u> | COMP1, cooperates with myogenic proteins in multicomponent complex. |
| <u>V\$COUP</u> | Repr. of RXR- mediated activ. & retinoic acid responses | <u>V\$COUP.01</u> | COUP antagonizes HNF- 4 by binding site competition or synergizes by direct protein - protein interaction with HNF-4 |
| <u>V\$CP2F</u> | CP2-erythrocyte Factor related to drosophila Elf1 | <u>V\$CP2.01</u> | CP2 |
| <u>V\$CREB</u> | Camp-Responsive Element Binding proteins | <u>V\$CREBP1.01</u> | cAMP-responsive element binding protein 1 |
| | | <u>V\$CREBP1CJUN.01</u> | CRE-binding protein 1/c- Jun heterodimer |
| | | <u>V\$CREB.01</u> | cAMP-responsive element binding protein |
| | | <u>V\$HLF.01</u> | hepatic leukemia factor |
| | | <u>V\$E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| | | <u>V\$CREB.02</u> | cAMP-responsive element binding protein |
| | | <u>V\$CREB.03</u> | cAMP-response element- binding protein |
| | | <u>V\$CREB.04</u> | cAMP-response element binding protein |
| | | <u>V\$CREBP1.02</u> | CRE-binding protein 1 |
| | | <u>V\$ATF.02</u> | ATF binding site |
| | | <u>V\$ATF.01</u> | activating transcription factor |
| | | <u>V\$TAXCREB.01</u> | Tax/CREB complex |
| | | <u>V\$TAXCREB.02</u> | Tax/CREB complex |

| Family | Family Information | Matrix Name | Information |
|----------------|--|-----------------------|--|
| | | <u>V\$VJUN.01</u> | v-Jun |
| <u>V\$E2FF</u> | E2F-myc activator/cell cycle regulator | <u>V\$E2F.02</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| | | <u>V\$E2F.03</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| | | <u>V\$E2F.01</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| | | | |
| <u>V\$E2TF</u> | papilloma virus E2 Transcriptional activator | <u>V\$E2.01</u> | BPV bovine papilloma virus regulator E2 |
| | | <u>V\$E2.02</u> | papilloma virus regulator E2 |
| <u>V\$EBOR</u> | E-BOx Related factors | <u>V\$DELTAEF1.01</u> | deltaEF1 |
| | | <u>V\$XBP1.01</u> | X-box-binding protein 1 |
| <u>V\$EBOX</u> | E-BOX binding factors | <u>V\$USF.02</u> | upstream stimulating factor |
| | | <u>V\$USF.03</u> | upstream stimulating factor |
| | | <u>V\$MYCMAX.03</u> | MYC-MAX binding sites |
| | | <u>V\$SREBP.03</u> | Sterol regulatory element binding protein |
| | | <u>V\$SREBP.02</u> | Sterol regulatory element binding protein |
| | | <u>V\$MYCMAX.02</u> | c-Myc/Max heterodimer |
| | | <u>V\$NMYC.01</u> | N-Myc |
| | | <u>V\$ATF6.01</u> | Member of b-zip family, induced by ER damage/stress |
| | | <u>V\$USF.01</u> | upstream stimulating factor |

| Family | Family Information | Matrix Name | Information |
|----------------|--|-----------------------|---|
| | | <u>V\$MYCMAX.01</u> | c-Myc/Max heterodimer |
| | | <u>V\$MAX.01</u> | Max |
| | | <u>V\$ARNT.01</u> | AhR nuclear translocator homodimers |
| | | <u>V\$SREBP.01</u> | Sterol regulatory element binding protein 1 and 2 |
| <u>V\$ECAT</u> | Enhancer-CcAaT binding factors | <u>V\$NFY.02</u> | nuclear factor Y (Y-box binding factor) |
| | | <u>V\$NFY.03</u> | nuclear factor Y (Y-box binding factor) |
| | | <u>V\$NFY.01</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$EGRF</u> | EGR/nerve growth Factor Induced protein C & rel. fact. | <u>V\$EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| | | <u>V\$EGR2.01</u> | Egr-2/Krox-20 early growth response gene product |
| | | <u>V\$EGR3.01</u> | early growth response gene 3 product |
| | | <u>V\$NGFIC.01</u> | nerve growth factor-induced protein C |
| | | <u>V\$WT1.01</u> | Wilms Tumor Suppressor |
| <u>V\$EKLF</u> | Erythroid krueppel like factor | <u>V\$EKLF.01</u> | Erythroid krueppel like factor (EKLF) |
| <u>V\$ETSF</u> | Human and murine ETS1 Factors | <u>V\$CETS1P54.01</u> | c-Ets-1(p54) |
| | | <u>V\$NRF2.01</u> | nuclear respiratory factor 2 |
| | | <u>V\$GABP.01</u> | GABP: GA binding protein |
| | | <u>V\$ELK1.02</u> | Elk-1 |

| Family | Family Information | Matrix Name | Information |
|----------------|-----------------------------------|-------------------|---|
| | | <u>V\$FLI.01</u> | ETS family member FLI |
| | | <u>V\$ETS2.01</u> | c-Ets-2 binding site |
| | | <u>V\$ETS1.01</u> | c-Ets-1 binding site |
| | | <u>V\$ELK1.01</u> | Elk-1 |
| | | <u>V\$PU1.01</u> | Pu.1 (Pu120) Ets-like transcription factor identified in lymphoid B-cells |
| <u>V\$EVI1</u> | EVI1-myeloid transforming protein | <u>V\$EVI1.06</u> | Ecotropic viral integration site 1 encoded factor |
| | | <u>V\$EVI1.02</u> | Ecotropic viral integration site 1 encoded factor |
| | | <u>V\$EVI1.03</u> | Ecotropic viral integration site 1 encoded factor |
| | | <u>V\$EVI1.05</u> | Ecotropic viral integration site 1 encoded factor |
| | | <u>V\$EVI1.04</u> | Ecotropic viral integration site 1 encoded factor |
| | | <u>V\$EVI1.01</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$FKHD</u> | Fork Head Domain factors | <u>V\$HFH1.01</u> | HNF-3/Fkh Homolog 1 |
| | | <u>V\$HFH2.01</u> | HNF-3/Fkh Homolog 2 |
| | | <u>V\$HFH3.01</u> | HNF-3/Fkh Homolog 3 (= Freac-6) |
| | | <u>V\$HFH8.01</u> | HNF-3/Fkh Homolog-8 |
| | | <u>V\$XFD1.01</u> | Xenopus fork head domain factor 1 |

| Family | Family Information | Matrix Name | Information |
|----------------|--|----------------------|---|
| | | <u>V\$XFD2.01</u> | Xenopus fork head domain factor 2 |
| | | <u>V\$XFD3.01</u> | Xenopus fork head domain factor 3 |
| | | <u>V\$HNF3B.01</u> | Hepatocyte Nuclear Factor 3beta |
| | | <u>V\$FREAC2.01</u> | Fork head Related Activator-2 |
| | | <u>V\$FREAC3.01</u> | Fork head Related Activator-3 |
| | | <u>V\$FREAC4.01</u> | Fork head Related Activator-4 |
| | | <u>V\$FREAC7.01</u> | Fork head Related Activator-7 |
| <u>V\$GATA</u> | GATA binding factors | <u>V\$LMO2COM.02</u> | complex of Lmo2 bound to Tal-1, E2A proteins, and GATA-1, half-site 2 |
| | | <u>V\$GATA1.04</u> | GATA-binding factor 1 |
| | | <u>V\$GATA1.05</u> | GATA-binding factor 1 |
| | | <u>V\$GATA2.01</u> | GATA-binding factor 2 |
| | | <u>V\$GATA2.02</u> | GATA-binding factor 2 |
| | | <u>V\$GATA3.01</u> | GATA-binding factor 3 |
| | | <u>V\$GATA3.02</u> | GATA-binding factor 3 |
| | | <u>V\$GATA.01</u> | GATA binding site (consensus) |
| | | <u>V\$GATA1.03</u> | GATA-binding factor 1 |
| | | <u>V\$GATA1.01</u> | GATA-binding factor 1 |
| | | <u>V\$GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$GFI1</u> | Growth Factor Independence-transcriptional | <u>V\$GFI1.01</u> | growth factor independence 1 zinc finger protein acts as |

| Family | Family Information | Matrix Name | Information |
|----------------|---|---------------------|--|
| | repressor | | transcriptional repressor |
| <u>V\$GKLF</u> | Gut-enriched Krueppel Like binding Factor | <u>V\$GKLF.01</u> | gut-enriched Krueppel-like factor |
| <u>V\$GREF</u> | Glucocorticoid responsive and related elements | <u>V\$GRE.01</u> | Glucocorticoid receptor, C2C2 zinc finger protein binds glucocorticoid dependent to GREs |
| | | <u>V\$ARE.01</u> | Androgene receptor binding site |
| | | <u>V\$PRE.01</u> | Progesterone receptor binding site |
| <u>V\$HAML</u> | Human Acute Myelogenous Leukemia factors | <u>V\$AML1.01</u> | runt-factor AML-1 |
| <u>V\$HEAT</u> | HEAT shock factors | <u>V\$HSF1.01</u> | heat shock factor 1 |
| <u>V\$HEN1</u> | E-box binding factor without transcript. activation | <u>V\$HEN1.01</u> | HEN1 |
| | | <u>V\$HEN1.02</u> | HEN1 |
| <u>V\$HMTB</u> | Human muscle-specific Mt binding site | <u>V\$MTBF.01</u> | muscle-specific Mt binding site |
| <u>V\$HNF1</u> | Hepatic Nuclear Factor 1 | <u>V\$HNF1.01</u> | hepatic nuclear factor 1 |
| | | <u>V\$HNF1.02</u> | Hepatic nuclear factor 1 |
| <u>V\$HNF4</u> | Hepatic Nuclear Factor 4 | <u>V\$HNF4.01</u> | Hepatic nuclear factor 4 |
| | | <u>V\$HNF4.02</u> | Hepatic nuclear factor 4 |
| <u>V\$HOMS</u> | Homeodomain subfamily S8 | <u>V\$S8.01</u> | Binding site for S8 type homeodomains |
| <u>V\$HOXF</u> | Factors with moderate activity to homeo domain consensus sequence | <u>V\$HOXA9.01</u> | Member of the vertebrate HOX - cluster of homeobox factors |
| | | <u>V\$HOX1-3.01</u> | Hox-1.3, vertebrate homeobox protein |
| <u>V\$IKRS</u> | Ikaros zinc finger family | <u>V\$LYE1.01</u> | LYE-1 (Ikaros 1), enriched in B and T |

| Family | Family Information | Matrix Name | Information |
|----------------|---|---------------------|---|
| | | | lymphocytes |
| | | <u>V\$IK2.01</u> | Ikaros 2, potential regulator of lymphocyte differentiation |
| | | <u>V\$IK1.01</u> | Ikaros 1, potential regulator of lymphocyte differentiation |
| | | <u>V\$IK3.01</u> | Ikaros 3, potential regulator of lymphocyte differentiation |
| <u>V\$IRFF</u> | Interferon Regulatory Factors | <u>V\$IRF1.01</u> | interferon regulatory factor 1 |
| | | <u>V\$IRF2.01</u> | interferon regulatory factor 2 |
| | | <u>V\$ISRE.01</u> | interferon-stimulated response element |
| <u>V\$LEFF</u> | LEF1/TCF | <u>V\$LEF1.01</u> | TCF/LEF-1, involved in the Wnt signal transduction pathway |
| <u>V\$LTUP</u> | Lentiviral Tata UPstream element | <u>V\$TAACC.01</u> | Lentiviral TATA upstream element |
| <u>V\$MEF2</u> | MEF2-myocyte-specific enhancer-binding factor | <u>V\$MEF2.05</u> | MEF2 |
| | | <u>V\$MEF2.01</u> | myogenic enhancer factor 2 |
| | | <u>V\$HMEF2.01</u> | myocyte enhancer factor |
| | | <u>V\$MMEF2.01</u> | myocyte enhancer factor |
| | | <u>V\$RSRFC4.01</u> | related to serum response factor, C4 |
| | | <u>V\$RSRFC4.02</u> | related to serum response factor, C4 |
| | | <u>V\$AMEF2.01</u> | myocyte enhancer factor |
| | | <u>V\$MEF2.02</u> | myogenic MADS factor MEF-2 |

| Family | Family Information | Matrix Name | Information |
|----------------|---|-------------------------|---|
| | | <u>V\$MEF2.03</u> | myogenic MADS factor MEF-2 |
| | | <u>V\$MEF2.04</u> | myogenic MADS factor MEF-2 |
| <u>V\$MEF3</u> | MEF3 BINDING SITES | <u>V\$MEF3.01</u> | MEF3 binding site, present in skeletal muscle-specific transcriptional enhancers |
| <u>V\$MEIS</u> | Homeodomain factor aberrantly expressed in myeloid leukemia | <u>V\$MEIS1.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$MINI</u> | Muscle INItiator | <u>V\$MUSCLE_INI.01</u> | Muscle Initiator Sequence |
| | | <u>V\$MUSCLE_INI.02</u> | Muscle Initiator Sequence |
| | | <u>V\$MUSCLE_INI.03</u> | Muscle Initiator Sequence |
| <u>V\$MOKF</u> | Mouse Krueppel like factor | <u>V\$MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 |
| <u>V\$MTF1</u> | Metal induced transcription factor | <u>V\$MTF-1.01</u> | Metal transcription factor 1, MRE |
| <u>V\$MYOD</u> | MYOblast Determining factor | <u>V\$MYOD.02</u> | myoblast determining factor |
| | | <u>V\$MYF5.01</u> | Myf5 myogenic bHLH protein |
| | | <u>V\$MYOD.01</u> | myoblast determination gene product |
| | | <u>V\$LMO2COM.01</u> | complex of Lmo2 bound to Tal-1, E2A proteins, and GATA-1, half-site 1 |
| | | <u>V\$E47.01</u> | MyoD/E47 and MyoD/E12 dimers |
| | | <u>V\$E47.02</u> | TAL1/E47 dimers |

| Family | Family Information | Matrix Name | Information |
|-----------------------|---------------------------------------|--------------------------------|--|
| <u>V\$MYOF</u> | MYOgenic Factors | <u>V\$NF1.01</u> | nuclear factor 1 |
| | | <u>V\$MYOGNF1.01</u> | myogenin / nuclear factor 1 or related factors |
| <u>V\$MYT1</u> | Xenopus MYT1 C2HC zinc finger protein | <u>V\$MYT1.02</u> | MyT1 zinc finger transcription factor involved in primary neurogenesis |
| | | <u>V\$MYT1.01</u> | MyT1 zinc finger transcription factor involved in primary neurogenesis |
| <u>V\$MZF1</u> | Myeloid Zinc Finger 1 factors | <u>V\$MZF1.01</u> | MZF1 |
| <u>V\$NFAT</u> | Nuclear Factor of Activated T-cells | <u>V\$NFAT.01</u> | Nuclear factor of activated T-cells |
| <u>V\$NFKB</u> | Nuclear Factor Kappa B/c-rel | <u>V\$CREL.01</u> | c-Rel |
| | | <u>V\$NFKAPPAB.01</u> | NF-kappaB |
| | | <u>V\$NFKAPPAB65.01</u> | NF-kappaB (p65) |
| | | <u>V\$NFKAPPAB50.01</u> | NF-kappaB (p50) |
| | | <u>V\$NFKAPPAB.02</u> | NF-kappaB |
| | | <u>V\$NFKAPPAB.03</u> | NF-kappaB |
| <u>V\$NKXH</u> | NKX - Homeodomain sites | <u>V\$NKX25.01</u> | homeo domain factor Nkx-2.5/Csx, tinman homolog, high affinity sites |
| | | <u>V\$NKX25.02</u> | homeo domain factor Nkx-2.5/Csx, tinman homolog low affinity sites |
| | | <u>V\$NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$NOLF</u> | Neuron-specific-OLFactory factor | <u>V\$OLF1.01</u> | olfactory neuron-specific factor |

| Family | Family Information | Matrix Name | Information |
|-----------------------|---|---------------------------|--|
| <u>V\$NRSF</u> | Neuron-Restrictive Silencer Factor | <u>V\$NRSF.01</u> | neuron-restrictive silencer factor |
| | | <u>V\$NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$OAZF</u> | Olfactory associated zinc finger protein | <u>V\$ROAZ.01</u> | Rat C2H2 Zn finger protein involved in olfactory neuronal differentiation |
| <u>V\$OCT1</u> | OCTamer binding protein | <u>V\$OCT1.02</u> | octamer-binding factor 1 |
| | | <u>V\$OCT1.06</u> | octamer-binding factor 1 |
| | | <u>V\$OCT.01</u> | Octamer binding site (OCT1/OCT2 consensus) |
| | | <u>V\$OCT1.05</u> | octamer-binding factor 1 |
| | | <u>V\$OCT1.04</u> | octamer-binding factor 1 |
| | | <u>V\$OCT1.03</u> | octamer-binding factor 1 |
| | | <u>V\$OCT1.01</u> | octamer-binding factor 1 |
| <u>V\$OCTB</u> | OCT6 Binding factors_astrocytes + glioblastoma cells | <u>V\$TST1.01</u> | POU-factor Tst-1/Oct-6 |
| <u>V\$OCTP</u> | OCT1 binding factor (POU-specific domain) | <u>V\$OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$P53F</u> | p53 tumor suppr.-neg. regulat. of the tumor suppr. Rb | <u>V\$P53.01</u> | tumor suppressor p53 |
| <u>V\$PAX1</u> | PAX-1 binding site | <u>V\$PAX1.01</u> | Pax1 paired domain protein, expressed in the developing vertebral column of mouse embryos |
| <u>V\$PAX3</u> | PAX-3 binding sites | <u>V\$PAX3.01</u> | Pax-3 paired domain protein, expressed in embryogenesis, mutations correlate to Waardenburg Syndrome |

| Family | Family Information | Matrix Name | Information |
|----------------|--|------------------------|---|
| <u>VSPAX4</u> | Heterogeneous PAX-4 binding sites | <u>VSPAX4.01</u> | Pax-4 paired domain protein, together with PAX-6 involved in pancreatic development |
| <u>VSPAX5</u> | PAX-5 / PAX-9 B-cell-specific activating protein | <u>VSPAX9.01</u> | zebrafish PAX9 binding sites |
| | | <u>VSPAX5.01</u> | B-cell-specific activating protein |
| | | <u>VSPAX5.02</u> | B-cell-specific activating protein |
| <u>VSPAX6</u> | Activ. involved in Iris development in the mouse eye | <u>VSPAX6.01</u> | Pax-6 paired domain protein |
| <u>VSPAX8</u> | PAX-2/5/8 binding sites | <u>VSPAX8.01</u> | PAX 2/5/8 binding site |
| <u>VSPBXF</u> | Homeo domain factor PBX-1 | <u>VSPBX1.01</u> | homeo domain factor Pbx-1 |
| <u>VSPCAT</u> | Promoter-CcAaT binding factors | <u>VSACAAT.01</u> | Avian C-type LTR CCAAT box |
| | | <u>VSCAAT.01</u> | cellular and viral CCAAT box |
| | | <u>V\$CLTR CAAT.01</u> | Mammalian C-type LTR CCAAT box |
| <u>VSPDX1</u> | Pancreatic and intestinal homeodomain transcr. factor | <u>VSPDX1.01</u> | Pdx1 (IDX1/IPF1) pancreatic and intestinal homeodomain TF |
| | | <u>V\$ISL1.01</u> | Pancreatic and intestinal lim-homeodomain factor |
| <u>VSPERO</u> | PEROxisome proliferator-activated receptor | <u>V\$PPARA.01</u> | PPAR/RXR heterodimers |
| <u>VSPIT1</u> | GHF-1 pituitary specific pou domain transcription factor | <u>VSPIT1.01</u> | Pit1, GHF-1 pituitary specific pou domain transcription factor |
| <u>V\$RARF</u> | Nuclear receptor for | <u>V\$RAR.01</u> | Retinoic acid receptor |

| Family | Family Information | Matrix Name | Information |
|----------------|--|----------------------|--|
| | retenoic acid | | member of nuclear receptors |
| | | <u>V\$RTR.01</u> | Retinoid receptor-related testis-associated receptor (GCNF/RTR) |
| <u>V\$RBIT</u> | Regulator of B-Cell IgH transcription | <u>V\$BRIGHT.01</u> | Bright, B cell regulator of IgH transcription |
| <u>V\$RBPJ</u> | RBPJ - kappa | <u>V\$RBPJK.01</u> | Mammalian transcriptional repressor RBP-Jkappa/CBF1 |
| <u>V\$REBV</u> | Epstein-Barr virus transcription factor R | <u>V\$EBVR.01</u> | Epstein-Barr virus transcription factor R |
| <u>V\$RORA</u> | Estrogen receptor and rar-Rel. Orphan Receptor Alpha | <u>V\$RORA1.01</u> | RAR-related orphan receptor alpha1 |
| | | <u>V\$RORA2.01</u> | RAR-related orphan receptor alpha2 |
| | | <u>V\$ER.01</u> | estrogen receptor |
| <u>V\$RREB</u> | Ras-REsponsive element Binding protein | <u>V\$RREB1.01</u> | Ras-responsive element binding protein 1 |
| <u>V\$RXRF</u> | RXR heterodimer binding sites | <u>V\$FXRE.01</u> | Farnesoid X - activated receptor (RXR/FXR dimer) |
| | | <u>V\$VDR_RXR.01</u> | VDR/RXR Vitamin D receptor RXR heterodimer site |
| | | <u>V\$VDR_RXR.02</u> | VDR/RXR Vitamin D receptor RXR heterodimer site |
| | | <u>V\$LXRE.01</u> | Nuclear receptor involved in the regulation lipid homeostasis |
| <u>V\$SATB</u> | Special AT-rich sequence binding protein | <u>V\$SATB1.01</u> | Special AT-rich sequence-binding protein 1, predominantly expressed in thymocytes, binds to matrix |

| Family | Family Information | Matrix Name | Information |
|----------------|---|--------------------|---|
| | | | attachment regions (MARs) |
| <u>V\$SEF1</u> | SEF1 protein in mouse Retrovirus SL3-3 | <u>V\$SEF1.01</u> | SEF1 binding site |
| <u>V\$SF1F</u> | Vertebrate steroidogenic factor | <u>V\$SF1.01</u> | SF1 steroidogenic factor 1 |
| <u>V\$SMAD</u> | Vertebrate SMAD family of transcription factors | <u>V\$SMAD3.01</u> | Smad3 transcription factor involved in TGF-beta signaling |
| | | <u>V\$SMAD4.01</u> | Smad4 transcription factor involved in TGF-beta signaling |
| | | <u>V\$FAST1.01</u> | FAST-1 SMAD interacting protein |
| <u>V\$SORY</u> | SOx/sRY-sex/testis determinig and related HMG Box factors | <u>V\$SOX5.01</u> | Sox-5 |
| | | <u>V\$SRY.01</u> | sex-determining region Y gene product |
| | | <u>V\$HMGIY.01</u> | HMGI(Y) high-mobility-group protein I (Y), architectural transcription factor organizing the framework of a nuclear protein-DNA transcriptional complex |
| | | <u>V\$SOX9.01</u> | SOX (SRY-related HMG box) |
| <u>V\$SP1F</u> | GC-Box factors_SP1/GC | <u>V\$SP1.01</u> | stimulating protein 1 SP1, ubiquitous zinc finger transcription factor |
| | | <u>V\$GC.01</u> | GC box elements |
| <u>V\$SRFF</u> | Serum Response element binding Factor | <u>V\$SRF.02</u> | serum response factor |
| | | <u>V\$SRF.03</u> | serum responsive factor |
| | | <u>V\$SRF.01</u> | serum response factor |

| Family | Family Information | Matrix Name | Information |
|-----------------------|--|---------------------------|---|
| <u>V\$STAT</u> | Signal Transducer and Activator of Transcription factors | <u>V\$STAT.01</u> | signal transducers and activators of transcription |
| | | <u>V\$STAT5.01</u> | STAT5: signal transducer and activator of transcription 5 |
| | | <u>V\$STAT6.01</u> | STAT6: signal transducer and activator of transcription 6 |
| | | <u>V\$STAT1.01</u> | signal transducer and activator of transcription 1 |
| | | <u>V\$STAT3.01</u> | signal transducer and activator of transcription 3 |
| <u>V\$T3RH</u> | Viral homolog of thyroid hormone receptor alpha1 (AEV vErbA) | <u>V\$T3R.01</u> | vErbA, viral homolog of thyroid hormone receptor alpha1 |
| <u>V\$TBPF</u> | Tata-Binding Protein Factor | <u>V\$TATA.02</u> | Mammalian C-type LTR TATA box |
| | | <u>V\$ATATA.01</u> | Avian C-type LTR TATA box |
| | | <u>V\$TATA.01</u> | cellular and viral TATA box elements |
| | | <u>V\$MTATA.01</u> | Muscle TATA box |
| <u>V\$TCFF</u> | TCF11 transcription Factor | <u>V\$TCF11.01</u> | TCF11/KCR-F1/Nrf1 homodimers |
| <u>V\$TEAF</u> | TEA/ATTS DNA binding domain factors | <u>V\$TEF1.01</u> | TEF-1 related muscle factor |
| <u>V\$TTFF</u> | Thyroid transcription factor-1 | <u>V\$TTF1.01</u> | Thyroid transcription factor-1 (TTF1) binding site |
| <u>V\$VBPF</u> | chicken Vitellogenin gene Binding Protein factor | <u>V\$VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |

| Family | Family Information | Matrix Name | Information |
|-----------------------|--|---------------------------|---|
| <u>V\$VMYB</u> | AMV-viral myb oncogene | <u>V\$VMYB.02</u> | v-Myb |
| | | <u>V\$VMYB.01</u> | v-Myb |
| <u>V\$WHZF</u> | Winged Helix and ZF5 binding sites | <u>V\$WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$XBBF</u> | X-box binding Factors | <u>V\$RFX1.01</u> | X-box binding protein RFX1 |
| | | <u>V\$RFX1.02</u> | X-box binding protein RFX1 |
| | | <u>V\$MIF1.01</u> | MIBP-1 / RFX1 complex |
| <u>V\$XSEC</u> | Xenopus SEleno Cystein t-RNA activating factor | <u>V\$STAF.02</u> | Se-Cys tRNA gene transcription activating factor |
| | | <u>V\$STAF.01</u> | Se-Cys tRNA gene transcription activating factor |
| <u>V\$YY1F</u> | activator/repressor binding to transcr. init. site | <u>V\$YY1.01</u> | Yin and Yang 1 |
| <u>V\$ZBPF</u> | Zinc binding protein factor | <u>V\$ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$ZFIA</u> | ZincFinger with InterAction domain factors | <u>V\$ZID.01</u> | zinc finger with interaction domain |

© Genomatix Software GmbH 1998-2002 - All rights reserved.

B. Changes from Family Library Version 2.4 to Version 3.0

Matrix Family Library Version 3.0 (Nov 2002) contains 452 weight matrices in 216 families

(Vertebrates: 314 matrices in 128 families)

5 **New weight matrices - Vertebrates**

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|--|--------------------|---|
| <u>V\$AP1F</u> | AP1 and related factors | <u>V\$BACH1.01</u> | BTB/POZ-bZIP transcription factor BACH1 forms heterodimers with the small Maf protein family |
| <u>V\$CIZF</u> | CAS interacting zinc finger protei | <u>V\$NMP4.01</u> | NMP4 (nuclear matrix protein 4) / CIZ (Cas-interacting zinc finger protein) |
| <u>V\$CREB</u> | Camp-Responsive Element Binding proteins | <u>V\$ATF6.02</u> | Activating transcription factor 6, member of b-zip family, induced by ER stress |
| <u>V\$E4FF</u> | Ubiquitous GLI - Krueppel like zinc finger involved in cell cycle regulation | <u>V\$E4F.01</u> | GLI-Krueppel-related transcription factor, regulator of adenovirus E4 promoter |
| <u>V\$GFI1</u> | Growth Factor Independence-transcriptional repressor | <u>V\$Gfi1B.01</u> | Growth factor independence 1 zinc finger protein Gfi-1B |
| <u>V\$GLIF</u> | GLI zinc finger family | <u>V\$GLI1.01</u> | Zinc finger transcription factor GLI1 |
| <u>V\$HAML</u> | Human Acute Myelogenous Leukemia factors | <u>V\$AML3.01</u> | Runt-related transcription factor 2 / CBFA1 (core-binding factor, runt domain, alpha subunit 1) |
| <u>V\$HESF</u> | Vertebrate homologues of enhancer of split complex | <u>V\$HES1.01</u> | Drosophila hairy and enhancer of split homologue 1 (HES-1) |
| <u>V\$HIF</u> | Hypoxia inducible factor, bHLH / PAS protein family | <u>V\$HIF1.01</u> | Hypoxia induced factor-1 (HIF-1) |
| | | <u>V\$HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|--|----------------------|--|
| <u>V\$HNF6</u> | Onecut Homeodomain factor HNF6 | <u>V\$HNF6.01</u> | Liver enriched Cut - Homeodomain transcription factor HNF6 (ONECUT) |
| <u>V\$HOXF</u> | Factors with moderate activity to homeo domain consensus sequence | <u>V\$CRX.01</u> | Cone-rod homeobox- containing transcription factor / otx-like homeobox gene |
| | | <u>V\$EN1.01</u> | Homeobox protein engrailed (en-1) |
| | | <u>V\$PTX1.01</u> | Pituitary Homeobox 1 (Ptx1) |
| <u>V\$IRFF</u> | Interferon Regulatory Factors | <u>V\$IRF3.01</u> | Interferon regulatory factor 3 (IRF-3) |
| | | <u>V\$IRF7.01</u> | Interferon regulatory factor 7 (IRF-7) |
| <u>V\$MAZF</u> | Myc associated zinc fingers | <u>V\$MAZ.01</u> | Myc associated zinc finger protein (MAZ) |
| | | <u>V\$MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$MEIS</u> | Homeodomain factor aberrantly expressed in myeloid leukemia | <u>V\$MEIS1.01</u> | Binding site for monomeric Meis1 homeodomain protein |
| <u>V\$MITF</u> | Microphthalmia transcription factor | <u>V\$MIT.01</u> | MIT (microphthalmia transcription factor) and TFE3 |
| <u>V\$MOKF</u> | Mouse Krueppel like factor | <u>V\$MOK2.02</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (human) |
| <u>V\$NEUR</u> | NeuroD, Beta2, HLH domain | <u>V\$NEUROD1.01</u> | DNA binding site for NEUROD1 (BETA-2 / E47 dimer) |
| <u>V\$NF1F</u> | Nuclear Factor 1 | <u>V\$NF1.02</u> | Nuclear factor 1 (CTF1) |
| <u>V\$NKXH</u> | NKX/DLX - Homeodomain sites | <u>V\$DLX1.01</u> | DLX-1, -2, and -5 binding sites |
| | | <u>V\$DLX3.01</u> | Distal-less 3 homeodomain transcription facto |
| | | <u>V\$HMX3.01</u> | H6 homeodomain HMX3/Nkx5.1 |

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|---|-------------------------|---|
| | | | transcription factor |
| | | <u>V\$MSX.01</u> | Homeodomain proteins MSX-1 and MSX-2 |
| | | <u>V\$MSX2.01</u> | Muscle segment homeo box 2, homologue of Drosophila (HOX 8) |
| <u>V\$NRLF</u> | Neural retina leucine zipper | <u>V\$NRL.01</u> | Neural retinal basic leucine zipper factor (bZIP) |
| <u>V\$PARF</u> | PAR/bZIP family | <u>V\$DBP.01</u> | Albumin D-box binding protein |
| | | <u>V\$PBX1_MEIS1.01</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$PBXC</u> | PBX1 - MEIS1 complexes | <u>V\$PBX1_MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| | | <u>V\$PBX1_MEIS1.03</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$PLZF</u> | C2H2 zinc finger protein PLZF | <u>V\$PLZF.01</u> | Promyelocytic leukemia zink finger (TF with nine Krueppel-like zink fingers) |
| | | | Halfsite of PXR (pregnane X receptor)/RXR resp. CAR (constitutive androstane receptor)/RXR heterodimer binding site |
| <u>V\$PXR</u> | Pregnane X receptor | <u>V\$PXR CAR.01</u> | |
| <u>V\$RORA</u> | v-ERB and rar-related Orphan Receptor Alpha | <u>V\$NBRE.01</u> | Monomers of the nur subfamily of nuclear receptors (nur77, nurr1, nor-1) |
| <u>V\$SF1F</u> | Vertebrate steroidogenic factor | <u>V\$FTF.01</u> | Alpha (1)-fetoprotein transcription factor (FTF), liver receptor homologue-1 (LHR-1) |
| <u>V\$SIXF</u> | Sine oculis (SIX) homeodomain factors | <u>V\$SIX3.01</u> | SIX3 / SIXdomain (SD) and Homeodomain (HD) transcription factor |
| <u>V\$TALE</u> | TALE Homeodomain class recognizing TG motives | <u>V\$TGIF.01</u> | TG-interacting factor belonging to TALE class of homeodomain factors |

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|----------------------------|------------------|---|
| <u>V\$ZF5F</u> | ZF5 POZ domain zinc finger | <u>V\$ZF5.01</u> | Zinc finger / POZ domain transcription factor |

Weight matrices renamed

- V\$MEIS1.01 renamed to V\$MEIS1_HOXA9.01

Weight matrices moved to other families

- V\$BEL1.01 moved from V\$AP1F to V\$BEL1
- 5 • V\$NF1.01 moved from V\$MYOF to V\$NF1
- V\$ER.01 moved from V\$RORA to V\$EREF
- V\$T3R.01 moved from V\$T3RH to V\$RORA
- V\$CLTR_CAAT.01 moved from V\$PCAT to V\$RCAT
- V\$FAST1.01 moved from V\$SMAD to V\$FAST

10 **Weight matrices removed**

- V\$MUSCLE_INI.03

C. Changes from Family Library Version 3.0 to Version 3.1

- 15 Matrix Family Library Version 3.1 contains 456 weight matrices in 216 families (Vertebrates: 318 matrices in 128 families)

New weight matrices - Vertebrates

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|--|----------------------|--|
| <u>V\$LEFF</u> | LEF1/TCF | <u>V\$LEF1.02</u> | TCF/LEF-1, involved in the Wnt signal transduction pathway |
| <u>V\$PAX2</u> | PAX-2 binding sites | <u>V\$PAX2.01</u> | Zebrafish PAX2 paired domain protein |
| <u>V\$PAX5</u> | PAX-5/PAX-9 B-cell-specific activating protein | <u>V\$PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$PAX6</u> | PAX-4/PAX-6 paired domain binding sites | <u>V\$PAX4_PD.01</u> | PAX4 paired domain binding site |
| | | <u>V\$PAX6.02</u> | PAX6 paired domain and homeodomain are required for binding to this site |
| <u>V\$ZBPF</u> | Zinc binding protein | <u>V\$ZF9.01</u> | Core promoter-binding |

| | | | |
|--|--------|--|---|
| | factor | | protein (CPBP) with 3 Krueppel-type zinc fingers |
|--|--------|--|---|

Weight matrices modified

- V\$AML1.01
- V\$AML3.01

5 Weight matrices moved to other families

- V\$ARNT.01 moved from V\$EBOX to V\$HIF1 (ARNT is a synonym for HIF1 B)

Weight matrices removed

- V\$SEF1.01
- 10 • V\$OCT1.03

Version 3.1.1 (April 2003)

Matrices V\$IRF3.01 and V\$IRF7.01 corrected.

Version 3.1.2 (June 2003)

Matrix V\$Gfi1B.01 corrected.

15

D. Changes from Family Library Version 3.1 to Version 3.3

Matrix Family Library Version 3.3 (August 2003) contains 485 weight matrices in 233 families

(Vertebrates: 326 matrices in 130 families)

5 New weight matrices - Vertebrates

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|---|--------------------------|---|
| <u>V\$EREF</u> | Estrogen Response Elements | <u>V\$ER.02</u> | Canonical palindromic estrogen response element (ERE) |
| <u>V\$SP1F</u> | GC-Box factors_SP1/GC | <u>V\$BTEB3.01</u> | Basic transcription element (BTE) binding protein, BTEB3, FKL-2 |
| <u>V\$CDEF</u> | Cell cycle regulators: Cell cycle dependent element | <u>V\$CDE.01</u> | Cell cycle-dependent element, CDF-1 binding site (CDE/CHR tandem elements regulate cell cycle dependent repression) |
| <u>V\$CHRF</u> | Cell cycle regulators: Cell cycle homology element | <u>V\$CHR.01</u> | Cell cycle gene homology region (CDE/CHR tandem elements regulate cell cycle dependent repression) |
| <u>V\$HIF</u> | Hypoxia inducible factor, bHLH / PAS protein family | <u>V\$CLOCK_BMAL1.01</u> | Binding site of Clock/BMAL1 heterodimer, NPAS2/BMAL1 heterodimer |
| <u>V\$FKHD</u> | Fork Head Domain factors | <u>V\$FKHRL1.01</u> | Fkh-domain factor FKHRL1 (FOXO) |
| <u>V\$P53F</u> | p53 tumor suppr.-neg. regulat. of the tumor suppr. Rb | <u>V\$P53.02</u> | Tumor suppressor p53 (5' half site) |
| | | <u>V\$P53.03</u> | Tumor suppressor p53 (3' half site) |

Weight matrices modified

- V\$GFI1.01

E. Changes from Family Library Version 3.3 to Version 4.0

Matrix Family Library Version 4.0 (November 2003) contains 535 weight matrices in 253 families

(Vertebrates: 339 matrices in 136 families)

5 New weight matrices - Vertebrates

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|--|---------------------|---|
| <u>V\$AARE</u> | AARE binding factors | <u>V\$AARE.01</u> | Amino acid response element, ATF4 binding site |
| <u>V\$AP1R</u> | MAF and AP1 related factors | <u>V\$BACH2.01</u> | Bach2 bound TRE |
| | | <u>V\$NFE2L2.01</u> | Nuclear factor (erythroid-derived 2)-like 2, NRF2 |
| <u>V\$CDXF</u> | Vertebrate caudal related homeodomain protein | <u>V\$CDX1.01</u> | Intestine specific homeodomain factor CDX-1 |
| <u>V\$DEAF</u> | Homolog to deformed epidermal autoregulatory factor-1 from D. melanogaster | <u>V\$NUDR.01</u> | NUDR (nuclear DEAF-1 related transcriptional regulator protein |
| <u>V\$ETSF</u> | Human and murine ETS1 factors | <u>V\$ELF2.01</u> | Ets - family member ELF-2 (NERF1a) |
| <u>V\$GABF</u> | GA-boxes | <u>V\$GAGA.01</u> | GAGA-Box |
| <u>V\$HNF1</u> | Hepatic Nuclear Factor 1 | <u>V\$HNF1.03</u> | Hepatic nuclear factor 1 |
| <u>V\$HOXF</u> | Factors with moderate activity to homeo domain consensus sequence | <u>V\$GSC.01</u> | Vertebrate bicoid-type homeodomain protein Goosecoid |
| <u>V\$LHXF</u> | Lim homeodomain factors | <u>V\$LHX3.01</u> | Homeodomain binding site in LIM/Homeodomain factor LHX3 |
| <u>V\$NKXH</u> | NKX/DLX - homeodomain sites | <u>V\$NKX32.01</u> | Homeodomain protein NKX3.2 (BAPX1, NKX3B, Bagpipe homolog) |
| <u>V\$RBPF</u> | RBPJ - kappa | <u>V\$RBPK.02</u> | Mammalian transcriptional repressor RBP-Jkappa/CBF1 |
| <u>V\$RP58</u> | RP58 (ZFP238) zinc finger protein | <u>V\$RP58.01</u> | Zinc finger protein RP58 (ZNF238), associated preferentially with heterochromatin |

Weight matrices modified

- V\$GRE.01

- V\$NFY.03

Weight matrices moved to other families

- V\$BACH1.01 moved from V\$AP1F to V\$AP1R
- V\$NFE2.01 moved from V\$AP1F to V\$AP1R
- 5 • V\$TCF11MAFG.01 moved from V\$AP1F to V\$AP1R
- V\$VMAF.01 moved from V\$AP1F to V\$AP1R

E. Changes from Family Library Version 4.0 to Version 4.1

- 10 Matrix Family Library Version 4.1 (February 2004) contains 564 weight matrices in 262 families
(Vertebrates: 356 matrices in 138 families)

New weight matrices - Vertebrates

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|---|------------------------|---|
| <u>V\$BNCF</u> | Basonuclein rDNA transcription factor (PolI) | <u>V\$BNC.01</u> | Basonuclin, cooperates with USF1 in rDNA PolI transcription) |
| <u>V\$CMYB</u> | C-myb, cellular transcriptional activator | <u>V\$CMYB.02</u> | c-Myb, important in hematopoiesis, cellular equivalent to avian myoblastosis virus oncogene v-myb |
| <u>V\$CP2F</u> | CP2-erythrocyte Factor related to drosophila Elf1 | <u>V\$CP2.02</u> | LBP-1c (leader-binding protein-1c), LSF (late SV40 factor), CP2, SEF (SAA3 enhancer factor) |
| <u>V\$EKLF</u> | Basic and erythroid Krueppel like factors | <u>V\$BKLF.01</u> | Basic krueppel-like factor (KLF3) |
| <u>V\$HAND</u> | bHLH transcription factor dimer of HAND2 and E12 | <u>V\$HAND2 E12.01</u> | Heterodimers of the bHLH transcription factors HAND2 (Thing2) and E12 |
| <u>V\$HIF</u> | Hypoxia inducible factor, bHLH / PAS protein family | <u>V\$DEC1.01</u> | Basic helix-loop-helix protein known as Dec1, Stra13 or Sharp2 |
| <u>V\$HNF6</u> | Onecut Homeodomain factor HNF6 | <u>V\$OC2.01</u> | CUT-homeodomain transcription factor Onecut-2 |
| <u>V\$HOXF</u> | Factors with moderate activity to homeo domain consensus sequence | <u>V\$OTX2.01</u> | Homeodomain transcription factor Otx2 (homolog of Drosophila orthodenticle) |

| Family | Family Information | Matrix Name | Matrix Information |
|----------------|-------------------------------|---------------------|--|
| | | <u>V\$GSH1.01</u> | Homeobox transcription factor Gsh-1 |
| <u>V\$IRFF</u> | Interferon Regulatory Factors | <u>V\$IRF4.01</u> | Interferon regulatory factor (IRF)-related protein (NF-EM5, PIP, LSIRF, ICSAT) |
| <u>V\$LHFX</u> | Lim homeodomain factors | <u>V\$LMX1B.01</u> | LIM-homeodomain transcription factor |
| <u>V\$MYT1</u> | MYT1 C2HC zinc finger protein | <u>V\$MYT1L.01</u> | Myelin transcription factor 1-like, neuronal C2HC zinc finger factor 1 |
| <u>V\$NEUR</u> | NeuroD, Beta2, HLH domain | <u>V\$NEUROG.01</u> | Neurogenin 1 and 3 (ngn1/3) binding sites |
| <u>V\$VMYB</u> | AMV-viral myb oncogene | <u>V\$VMYB.03</u> | v-Myb, viral myb variant from transformed BM2 cells |
| | | <u>V\$VMYB.04</u> | v-Myb, AMV v-myb |
| | | <u>V\$VMYB.05</u> | v-Myb, variant of AMV v-myb |
| <u>V\$ZBPF</u> | Zinc binding protein factor | <u>V\$ZNF202.01</u> | Transcriptional repressor, binds to elements found predominantly in genes that participate in lipid metabolism |

Weight matrices modified

- V\$CMYB.01
- V\$PTX1.01

Copyright © Genomatix Software GmbH 1998-2004 - All rights reserved

5

Example 6

Summary of Design for Particular Selectable Genes

TF binding sites and search parameters

Each TF binding site ("matrix") belongs to a matrix family that groups functionally similar matrices together, eliminating redundant matches by MatInspector professional (the search program). Searches were limited to vertebrate TF binding sites. Searches were performed by matrix family, i.e., the results show only the best match from a family for each site. MatInspector

10

default parameters were used for the core and matrix similarity values (core similarity = 0.75, matrix similarity = optimized).

Table 18

Gene Designations

5

A. Synthetic hygromycin gene

| Sequence | Description | Matrix Library |
|----------|--|--------------------|
| hyg | from pcDNA3.1/Hygro | Not applicable |
| hhyg | humanized ORF | Not applicable |
| hhyg-1 | First removal of undesired sequence matches | Ver 3.1.2 Jun 2003 |
| hhyg-2 | Second removal of undesired sequence matches | Ver 3.1.2 Jun 2003 |
| hhyg-3 | Third removal of undesired sequence matches | Ver 3.1.2 Jun 2003 |
| hHygro | Changes to ORF and add linker | Ver 3.3 Aug 2003 |
| hhyg-4 | Fourth removal of undesired sequence matches | Ver 3.3 Aug 2003 |

B. Synthetic neomycin gene

| Sequence | Description | Matrix Library |
|----------|--|--------------------|
| neo | from pCI-neo or psiSTRIKE neo | Not applicable |
| hneo | humanized ORF | Not applicable |
| hneo-1 | First removal of undesired sequence matches | Ver 3.1.2 Jun 2003 |
| hneo-2 | Second removal of undesired sequence matches | Ver 3.1.2 Jun 2003 |
| hneo-3 | Third removal of undesired sequence matches | Ver 3.1.2 Jun 2003 |
| hneo-4 | Changed 5' and 3' flanking regions/cloning sites | Ver 4.1 Feb 2004 |
| hneo-5 | Fourth removal of undesired sequence matches | Ver 4.1 Feb 2004 |

C. Synthetic puromycin gene

| Sequence | Description | Matrix Library |
|----------|--|------------------|
| puro | from psiSTRIKE puromycin | Not applicable |
| hpuro | humanized ORF | Not applicable |
| hpuro-1 | First removal of undesired sequence matches | Ver 4.1 Feb 2004 |
| hpuro-2 | Second removal of undesired sequence matches | Ver 4.1 Feb 2004 |

Note: the above sequence names designate the ORF only (except for Hhygro which includes flanking sequences). Addition of "F" to the sequence name indicates the presence of up- and down-stream flanking sequences. Additional letters (e.g., "B") indicate changes were made only to the flanking regions

Table 19

Sequences in Synthetic Hygromycin GenesTFBS in hhyg

Before removal of TFBS from hhyg (94 matches)

| Family/matrix | Further Information |
|-----------------------|---|
| V\$PCAT/CAAT.01 | cellular and viral CCAAT box |
| V\$MINI/MUSCLE_INI.02 | Muscle Initiator Sequence |
| V\$MINI/MUSCLE_INI.01 | Muscle Initiator Sequence |
| V\$SETSF/PU1.01 | Pu.1 (Pu120) Ets-like transcription factor identified in lymphoid B-cells |
| V\$AHRR/AHRARNT.02 | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| V\$EGRF/EGR3.01 | early growth response gene 3 product |
| V\$AP4R/AP4.01 | Activator protein 4 |
| V\$EGRF/NGFIC.01 | Nerve growth factor-induced protein C |
| V\$MAZF/MAZ.01 | Myc associated zinc finger protein (MAZ) |
| V\$ZBPF/ZF9.01 | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| V\$CREB/ATF6.02 | Activating transcription factor 6, member of b-zip family, induced by ER stress |
| V\$EGRF/EGR3.01 | early growth response gene 3 product |
| V\$ZBPF/ZF9.01 | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| V\$HIF/HIF1.02 | Hypoxia inducible factor, bHLH / PAS |

| Family/matrix | Further Information |
|-----------------------------|---|
| | protein family |
| <u>V\$E2FF/E2F.01</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| <u>V\$AP4R/AP4.01</u> | Activator protein 4 |
| <u>V\$HEN1/HEN1.02</u> | HEN1 |
| <u>V\$MYOD/E47.01</u> | MyoD/E47 and MyoD/E12 dimers |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$MOKF/MOK2.02</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (human) |
| <u>V\$SP1F/GC.01</u> | GC box elements |
| <u>V\$NRSE/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$RORA/RORA2.01</u> | RAR-related orphan receptor alpha2 |
| <u>V\$ZBPZ/ZF9.01</u> | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$AHRR/AHRARNT.02</u> | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$AP1F/TCF11MAFG.01</u> | TCF11/MafG heterodimers, binding to subclass of AP1 sites |
| <u>V\$EKLF/EKLF.01</u> | Erythroid krueppel like factor (EKLF) |
| <u>V\$NRSE/NRSE.01</u> | Neuron-restrictive silencer factor |
| <u>V\$NRSE/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$EBOX/MYCMAX.03</u> | MYC-MAX binding sites |
| <u>V\$RXRF/FXRE.01</u> | Farnesoid X - activated receptor (RXR/FXR dimer) |
| <u>V\$AHRR/AHRARNT.02</u> | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$WHZF/WHN.01</u> | Winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$EGRF/EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| <u>V\$SMAD/SMAD3.01</u> | Smad3 transcription factor involved in TGF-beta signaling |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$MYOD/MYOD.02</u> | Myoblast determining factor |

| Family/matrix | Further Information |
|------------------------------|---|
| <u>V\$E4FF/E4F.01</u> | GLI-Krueppel-related transcription factor, regulator of adenovirus E4 promoter |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$EGRF/EGR2.01</u> | Egr-2/Krox-20 early growth response gene product |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$EBOX/USF.02</u> | Upstream stimulating factor |
| <u>V\$HIF/ARNT.01</u> | AhR nuclear translocator homodimers |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$BEL1/BEL1.01</u> | Bel-1 similar region (defined in Lentivirus LTRs) |
| <u>V\$NRSE/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$MYOD/MYOD.01</u> | Myoblast determination gene product |
| <u>V\$NEUR/NEUROD1.01</u> | DNA binding site for NEUROD1 (BETA-2 / E47 dimer) |
| <u>V\$AHR/ARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$HIF/ARNT.01</u> | AhR nuclear translocator homodimers |
| <u>V\$MYB/VMYB.02</u> | v-Myb |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$PBX/PBX1 MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$MYO/MYOGNF1.01</u> | Myogenin / nuclear factor 1 or related factors |
| <u>V\$SRF/SRF.03</u> | Serum responsive factor |
| <u>V\$CP2/CP2.01</u> | CP2 |
| <u>V\$OAZ/ROAZ.01</u> | Rat C2H2 Zn finger protein involved in olfactory neuronal differentiation |
| <u>V\$AHR/ARNT.01</u> | Aryl hydrocarbon / dioxin receptor |
| <u>V\$MINI/MUSCLE INI.01</u> | Muscle Initiator Sequence |
| <u>V\$PAX5/PAX5.02</u> | B-cell-specific activating protein |

| Family/matrix | Further Information |
|------------------------------|---|
| <u>V\$ZBPF/ZF9.01</u> | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$EGRF/NGFIC.01</u> | Nerve growth factor-induced protein C |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$AP4R/AP4.02</u> | Activator protein 4 |
| <u>V\$XBBF/MIF1.01</u> | MIBP-1 / RFX1 complex |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$WHZF/WHN.01</u> | Winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$WHZF/WHN.01</u> | Winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$PAX5/PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$PAX5/PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$ZBPF/ZF9.01</u> | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| <u>V\$CP2F/CP2.01</u> | CP2 |
| <u>V\$MINI/MUSCLE INI.02</u> | Muscle Initiator Sequence |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$AHRR/AHRARNT.02</u> | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$MINI/MUSCLE INI.02</u> | Muscle Initiator Sequence |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$SP1F/SP1.01</u> | stimulating protein 1 SP1, ubiquitous zinc finger transcription factor |
| <u>V\$ZBPF/ZF9.01</u> | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| <u>V\$EGRF/EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| <u>V\$EGRF/WT1.01</u> | Wilms Tumor Suppressor |

| Family/matrix | Further Information |
|----------------------|--|
| V\$SP1F/SP1.01 | stimulating protein 1 SP1, ubiquitous zinc finger transcription factor |
| V\$RCAT/CLTR CAAT.01 | Mammalian C-type LTR CCAAT box |
| V\$ZBPF/ZF9.01 | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| V\$EGRF/WT1.01 | Wilms Tumor Suppressor |
| V\$EGRF/WT1.01 | Wilms Tumor Suppressor |
| V\$NF1F/NF1.01 | Nuclear factor 1 |
| V\$PDX1/PDX1.01 | Pdx1 (IDX1/IPF1) pancreatic and intestinal homeodomain TF |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hhyg3

After removal of TFBS from hhyg2 (3 matches)

5

| Family/matrix | Further Information |
|-----------------------|------------------------------------|
| V\$MINI/MUSCLE INI.02 | Muscle Initiator Sequence |
| V\$PAX5/PAX5.02 | B-cell-specific activating protein |
| V\$VMYB/VMYB.02 | v-Myb |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hHygro

10 Before removal of TFBS from hHygro (5 matches, excluding linker)

| Family/matrix | Further Information |
|-----------------------|---|
| V\$MINI/MUSCLE INI.02 | Muscle Initiator Sequence |
| V\$PAX5/PAX5.02 | B-cell-specific activating protein |
| V\$AREB/AREB6.04 | AREB6 (Atp1a1 regulatory element binding factor 6) |
| V\$VMYB/VMYB.02 | v-Myb |
| V\$CDEF/CDE.01 | Cell cycle-dependent element, CDF-1 binding site (CDE/CHR tandem elements regulate cell cycle dependent repression) |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hhyg4

After removal of TFBS from hHygro (4 matches)

| Family/matrix | Further Information |
|------------------------------|--|
| <u>V\$MINI/MUSCLE_INI.02</u> | Muscle Initiator Sequence |
| <u>V\$PAX5/PAX5.02</u> | B-cell-specific activating protein |
| <u>V\$AREB/AREB6.04</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| <u>V\$VMYB/VMYB.02</u> | v-Myb |

**matches are listed in order of occurrence in the corresponding sequence

5

Table 20

Sequences in Synthetic Neomycin Genes10 TFBS in hneo

Before removal of TFBS from hneo (69 matches)

| Family/matrix | Further Information |
|------------------------------|---|
| <u>V\$PCAT/CAAT.01</u> | cellular and viral CCAAT box |
| <u>V\$ZFIA/ZID.01</u> | Zinc finger with interaction domain |
| <u>V\$AP1F/TCF11MAFG.01</u> | TCF11/MafG heterodimers, binding to subclass of AP1 sites |
| <u>V\$MINI/MUSCLE_INI.01</u> | Muscle Initiator Sequence |
| <u>V\$AHRR/AHRARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$SP1F/GC.01</u> | GC box elements |
| <u>V\$MINI/MUSCLE_INI.02</u> | Muscle Initiator Sequence |
| <u>V\$CP2F/CP2.01</u> | CP2 |
| <u>V\$WHZF/WHN.01</u> | Winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$PAX5/PAX5.02</u> | B-cell-specific activating protein |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$ZBPF/ZF9.01</u> | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| <u>V\$ZBPF/ZF9.01</u> | Core promoter-binding protein (CPBP) |

| Family/matrix | Further Information |
|---------------------------|---|
| | with 3 Krueppel-type zinc fingers |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$AHRR/AHRARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$NRSE/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$CREB/ATF6.02</u> | Activating transcription factor 6, member of b-zip family, induced by ER stress |
| <u>V\$RXR/VDR_RXR.01</u> | VDR/RXR Vitamin D receptor RXR heterodimer site |
| <u>V\$PCAT/CAAT.01</u> | cellular and viral CCAAT box |
| <u>V\$NRSE/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$P53/P53.01</u> | Tumor suppressor p53 |
| <u>V\$NEUR/NEUROD1.01</u> | DNA binding site for NEUROD1 (BETA-2 / E47 dimer) |
| <u>V\$EBOX/USF.03</u> | Upstream stimulating factor |
| <u>V\$MYOD/MYOD.02</u> | Myoblast determining factor |
| <u>V\$NRSE/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$WHZF/WHN.01</u> | Winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$EBOX/MYCMAX.03</u> | MYC-MAX binding sites |
| <u>V\$HESF/HES1.01</u> | Drosophila hairy and enhancer of split homologue 1 (HES-1) |
| <u>V\$NEUR/NEUROD1.01</u> | DNA binding site for NEUROD1 (BETA-2 / E47 dimer) |
| <u>V\$MYOD/MYOD.02</u> | Myoblast determining factor |
| <u>V\$REBV/EBVR.01</u> | Epstein-Barr virus transcription factor R |
| <u>V\$PAX5/PAX5.02</u> | B-cell-specific activating protein |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$EGRF/WT1.01</u> | Wilms Tumor Suppressor |
| <u>V\$EGRF/WT1.01</u> | Wilms Tumor Suppressor |
| <u>V\$ZBP/ZF9.01</u> | Core promoter-binding protein (CPBP) |

| Family/matrix | Further Information |
|------------------------------|--|
| | with 3 Krueppel-type zinc fingers |
| <u>V\$MINI/MUSCLE_INI.01</u> | Muscle Initiator Sequence |
| <u>V\$NRSF/NRSF.01</u> | Neuron-restrictive silencer factor |
| <u>USPflMI/PflMI</u> | RE II-IP |
| <u>V\$NRSF/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$MOKF/MOK2.02</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (human) |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$AP1F/AP1FJ.01</u> | Activator protein 1 |
| <u>V\$PAX5/PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$WHZF/WHN.01</u> | Winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$PAX6/PAX4_PD.01</u> | PAX4 paired domain binding site |
| <u>V\$VMYB/VMYB.02</u> | v-Myb |
| <u>V\$BEL1/BEL1.01</u> | Bel-1 similar region (defined in Lentivirus LTRs) |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$EGRF/EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$NRSF/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$ETSF/ETS1.01</u> | c-Ets-1 binding site |
| <u>V\$NRSF/NRSF.01</u> | Neuron-restrictive silencer factor |
| <u>V\$SP1F/SP1.01</u> | stimulating protein 1 SP1, ubiquitous zinc finger transcription factor |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$PAX5/PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$GREF/ARE.01</u> | Androgene receptor binding site |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$CLOX/CDP.01</u> | cut-like homeodomain protein |

****matches are listed in order of occurrence in the corresponding sequence**

TFBS in hneo3

After removal of TFBS from hneo2 = before removal of TFBS from hneo3 (0 matches)

TFBS in hneo4

After removal of TFBS from hneo3 = before removal of TFBS from hneo4 (7 matches)

| Family/matrix | Further Information |
|--------------------------|---|
| <u>V\$PAX5/PAX9.01</u> | Zebrafish PAX9 binding sites |
| <u>V\$AARF/AARE.01</u> | Amino acid response element, ATF4 binding site |
| <u>V\$P53F/P53.02</u> | Tumor suppressor p53 (5' half site) |
| <u>V\$AP1R/BACH2.01</u> | Bach2 bound TRE |
| <u>V\$NEUR/NEUROG.01</u> | Neurogenin 1 and 3 (ngn1/3) binding sites |
| <u>V\$CMYB/CMYB.01</u> | c-Myb, important in hematopoiesis, cellular equivalent to avian myoblastosis virus oncogene v-myb |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |

****matches are listed in order of occurrence in the corresponding sequence**

TFBS in hneo5

After removal of TFBS from hneo4 (0 matches)

Table 21

Sequences in Synthetic Puromycin Genes

TFBS matches in hpuro

Before removal of TFBS from hpuro (68 matches)

| Family/matrix | Further Information |
|------------------------|---|
| <u>V\$CDEF/CDE.01</u> | Cell cycle-dependent element, CDF-1 binding site (CDE/CHR tandem elements regulate cell cycle dependent repression) |
| <u>V\$PAX3/PAX3.01</u> | Pax-3 paired domain protein, expressed in embryogenesis, mutations correlate to Waardenburg |

| Family/matrix | Further Information |
|---------------------------|---|
| | Syndrome |
| <u>V\$CREB/ATF6.02</u> | Activating transcription factor 6, member of b-zip family, induced by ER stress |
| <u>V\$EBOR/XBP1.01</u> | X-box-binding protein 1 |
| <u>V\$P53F/P53.03</u> | Tumor suppressor p53 (3' half site) |
| <u>V\$HESF/HES1.01</u> | Drosophila hairy and enhancer of split homologue 1 (HES-1) |
| <u>V\$MTF1/MTF-1.01</u> | Metal transcription factor 1, MRE |
| <u>V\$EKLF/EKLF.01</u> | Erythroid krueppel like factor (EKLF) |
| <u>V\$EGRF/EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$CMYB/CMYB.01</u> | c-Myb, important in hematopoiesis, cellular equivalent to avian myoblastosis virus oncogene v-myb |
| <u>V\$AHRR/AHRARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$EBOX/MYCMAX.03</u> | MYC-MAX binding sites |
| <u>V\$RORA/RORA2.01</u> | RAR-related orphan receptor alpha2 |
| <u>V\$EBOX/MYCMAX.03</u> | MYC-MAX binding sites |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$EGRF/WT1.01</u> | Wilms Tumor Suppressor |
| <u>V\$HAML/AML3.01</u> | Runt-related transcription factor 2 / CBFA1 (core-binding factor, runt domain, alpha subunit 1) |
| <u>V\$PAX5/PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / |

| Family/matrix | Further Information |
|--------------------------------|--|
| | PAS protein family |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$OAZF/ROAZ.01</u> | Rat C2H2 Zn finger protein involved in olfactory neuronal differentiation |
| <u>V\$GABF/GAGA.01</u> | GAGA-Box |
| <u>V\$EBOX/MYCMAX.03</u> | MYC-MAX binding sites |
| <u>V\$MYOD/MYF5.01</u> | Myf5 myogenic bHLH protein |
| <u>V\$AP4R/TAL1BETA/E47.01</u> | Tal-1beta/E47 heterodimer |
| <u>V\$NEUR/NEUROG.01</u> | Neurogenin 1 and 3 (ngn1/3) binding sites |
| <u>V\$HAND/HAND2 E12.01</u> | Heterodimers of the bHLH transcription factors HAND2 (Thing2) and E12 |
| <u>V\$MAZF/MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$ZBPF/ZNF202.01</u> | Transcriptional repressor, binds to elements found predominantly in genes that participate in lipid metabolism |
| <u>V\$SP1F/SP1.01</u> | Stimulating protein 1 SP1, ubiquitous zinc finger transcription factor |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$RREB/RREB1.01</u> | Ras-responsive element binding protein 1 |
| <u>V\$XBBF/MIF1.01</u> | MIBP-1 / RFX1 complex |
| <u>V\$CREB/TAXCREB.01</u> | Tax/CREB complex |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$NRSF/NRSE.01</u> | Neural-restrictive-silencer-element |
| <u>V\$MINI/MUSCLE INI.02</u> | Muscle Initiator Sequence |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$DEAF/NUDR.01</u> | NUDR (nuclear DEAF-1 related |

| Family/matrix | Further Information |
|---------------------------|--|
| | transcriptional regulator protein) |
| <u>V\$AHRR/AHRARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$EGRF/EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$ETSF/ETS1.01</u> | c-Ets-1 binding site |
| <u>V\$STAT/STAT1.01</u> | Signal transducer and activator of transcription 1 |
| <u>V\$BCL6/BCL6.01</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$CREB/ATF6.02</u> | Activating transcription factor 6, member of b-zip family, induced by ER stress |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$EBOR/XBP1.01</u> | X-box-binding protein 1 |
| <u>V\$DEAF/NUDR.01</u> | NUDR (nuclear DEAF-1 related transcriptional regulator protein) |
| <u>V\$RXRF/VDR_RXR.01</u> | VDR/RXR Vitamin D receptor RXR heterodimer site |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$REBV/EBVR.01</u> | Epstein-Barr virus transcription factor R |
| <u>V\$ZBPF/ZF9.01</u> | Core promoter-binding protein (CPBP) with 3 Krueppel-type zinc fingers |
| <u>V\$MYOD/LMO2COM.01</u> | Complex of Lmo2 bound to Tal-1, |

| Family/matrix | Further Information |
|-------------------------|--|
| | E2A proteins, and GATA-1, half-site 1 |
| <u>V\$AREB/AREB6.03</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| <u>V\$RXRF/FXRE.01</u> | Farnesoid X - activated receptor (RXR/FXR dimer) |
| <u>V\$AHRR/AHR.01</u> | Aryl hydrocarbon / dioxin receptor |

**matches are listed in order of occurrence in the corresponding sequence

TFBS matches in hpuro1

After removal of TFBS from hpuro = before removal of TFBS from hpuro1
(4 matches)

5

| Family/matrix | Further Information |
|--------------------------|---|
| <u>V\$NEUR/NEUROG.01</u> | Neurogenin 1 and 3 (ngn1/3) binding sites |
| <u>V\$PAX5/PAX5.02</u> | B-cell-specific activating protein |
| <u>V\$REBV/EBVR.01</u> | Epstein-Barr virus transcription factor R |
| <u>V\$AHRR/AHR.01</u> | Aryl hydrocarbon / dioxin receptor |

**matches are listed in order of occurrence in the corresponding sequence

TFBS matches in hpuro2

After removal of TFBS from hpuro1 (2 matches)

10

| Family/matrix | Further Information |
|--------------------------|--|
| <u>V\$NEUR/NEUROG.01</u> | Neurogenin 1 and 3 (ngn1/3) binding sites |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |

**matches are listed in order of occurrence in the corresponding sequence

Example 7

15

Summary of Design of Synthetic Firefly Luciferase Genes

TF binding sites and search parameters

The TF binding sites are from the TF binding site library ("Matrix Family Library") that is part of the GEMS Launcher package. Each TF binding site ("matrix") belongs to a matrix family that groups functionally similar matrices

together, eliminating redundant matches by MatInspector professional (the search program). Searches were limited to vertebrate TF binding sites. Searches were performed by matrix family, i.e. the results show only the best match from a family for each site. MatInspector default parameters were used for the core and matrix similarity values (core similarity = 0.75, matrix similarity = optimized).

Table 22

Luc Gene Designations

10 Synthetic luc gene (versions A and B)

| Sequence* | Description | Matrix Library |
|--------------------------------------|---|--------------------|
| Luc | wild-type gene | (not applicable) |
| luc+ | improved gene from Promega's pGL3 vectors | (not applicable) |
| hluc+ | Improved gene form Promega's pGL3(R2.1)-Basic | (not applicable) |
| <i>Codon optimization strategy A</i> | | |
| hluc+ver2A1 | codon optimized luc+ (strategy A) | Ver 3.0 Nov 2002 |
| hluc+ver2A2 | First removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2A3 | Second removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2A4 | Third removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2A5 | Fourth removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2A6 | Fifth removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2A7 | Sixth removal of undesired sequence matches | Ver 3.1.1 Apr 2003 |
| hluc+ver2A8 | Removal of <i>Bgl</i> I (RE) site | Ver 3.1.1 Apr 2003 |
| <i>Codon optimization strategy B</i> | | |
| hluc+ver2B1 | codon optimized luc+ (strategy B) | Ver 3.0 Nov 2002 |
| hluc+ver2B2 | First removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2B3 | Second removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2B4 | Third removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2B5 | Fourth removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2B6 | Fifth removal of undesired sequence matches | Ver 3.0 Nov 2002 |
| hluc+ver2B7 | Sixth removal of undesired sequence matches | Ver 3.1.1 Apr 2003 |
| hluc+ver2B8 | Removal of <i>Sma</i> I (RE), Ptx1 (TF) sites | Ver 3.1.1 Apr 2003 |
| hluc+ver2B9 | Removal of additional CpG sequences | Ver 3.1.1 Apr 2003 |

| Sequence* | Description | Matrix Library |
|--------------|-----------------------------------|--------------------|
| hluc+ver2B10 | Removal of <i>Bgl</i> I (RE) site | Ver 3.1.1 Apr 2003 |

* the sequence names designate open reading frames; RE = restriction enzyme recognition sequence

Table 23

5

Sequences in Synthetic Luc Genes (version A)

TFBS in hluc+ver2A1

10

Before removal of TFBS from hluc+ver2A1 (110 matches)

| Family/matrix | Further Information |
|------------------------------|---|
| <u>V\$MINI/MUSCLE_INI.02</u> | Muscle Initiator Sequence |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$GREF/PRE.01</u> | Progesterone receptor binding site |
| <u>V\$MAZF/MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$SP1F/SP1.01</u> | stimulating protein 1 SP1, ubiquitous zinc finger transcription factor |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$SF1F/SF1.01</u> | SF1 steroidogenic factor 1 |
| <u>V\$EGRF/NGFIC.01</u> | Nerve growth factor-induced protein C |
| <u>V\$MINI/MUSCLE_INI.01</u> | Muscle Initiator Sequence |
| <u>V\$EGRF/EGR2.01</u> | Egr-2/Krox-20 early growth response gene product |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$HESF/HES1.01</u> | Drosophila hairy and enhancer of split homologue 1 (HES-1) |
| <u>V\$NRSF/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$PAX5/PAX5.02</u> | B-cell-specific activating protein |
| <u>V\$HAML/AML3.01</u> | Runt-related transcription factor 2 / CBFA1 (core-binding factor, runt domain, alpha subunit 1) |
| <u>V\$GREF/PRE.01</u> | Progesterone receptor binding site |
| <u>V\$P53F/P53.01</u> | tumor suppressor p53 |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$NF1F/NF1.01</u> | Nuclear factor 1 |

| Family/matrix | Further Information |
|------------------------------|---|
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$REB V/EBVR.01</u> | Epstein-Barr virus transcription factor R |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$PBXC/PBX1_MEIS1.01</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$XSEC/STAF.01</u> | Se-Cys tRNA gene transcription activating factor |
| <u>V\$COMP/COMP1.01</u> | COMP1, cooperates with myogenic proteins in multicomponent complex |
| <u>V\$MYOF/MYOGNF1.01</u> | Myogenin / nuclear factor 1 or related factors |
| <u>V\$NEUR/NEUROD1.01</u> | DNA binding site for NEUROD1 (BETA-2 / E47 dimer) |
| <u>V\$MYOD/MYOD.02</u> | myoblast determining factor |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$EVI1/EVI1.02</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$SMAD/SMAD4.01</u> | Smad4 transcription factor involved in TGF-beta signaling |
| <u>V\$MYOD/MYF5.01</u> | Myf5 myogenic bHLH protein |
| <u>V\$HESF/HES1.01</u> | Drosophila hairy and enhancer of split homologue 1 (HES-1) |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$SP1F/GC.01</u> | GC box elements |
| <u>V\$MAZF/MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$RREB/RREB1.01</u> | Ras-responsive element binding protein 1 |
| <u>V\$AHRR/AHRARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$YY1F/YY1.01</u> | Yin and Yang 1 |
| <u>V\$ETSF/GABP.01</u> | GABP: GA binding protein |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$ETSF/ELK1.02</u> | Elk-1 |
| <u>V\$EBOX/MYCMAX.03</u> | MYC-MAX binding sites |
| <u>V\$E4FF/E4F.01</u> | GLI-Krueppel-related transcription factor, regulator of adenovirus E4 promoter |

| Family/matrix | Further Information |
|------------------------------|--|
| <u>V\$XBBF/RFX1.01</u> | X-box binding protein RFX1 |
| <u>V\$EVI1/EVI1.06</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$NF1F/NF1.01</u> | Nuclear factor 1 |
| <u>V\$PBXC/PBX1_MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$HESF/HES1.01</u> | Drosophila hairy and enhancer of split homologue 1 (HES-1) |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$ETSF/GABP.01</u> | GABP: GA binding protein |
| <u>V\$MYOD/MYOD.02</u> | myoblast determining factor |
| <u>V\$XSEC/STAF.01</u> | Se-Cys tRNA gene transcription activating factor |
| <u>V\$OAZF/ROAZ.01</u> | Rat C2H2 Zn finger protein involved in olfactory neuronal differentiation |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$PAX3/PAX3.01</u> | Pax-3 paired domain protein, expressed in embryogenesis, mutations correlate to Waardenburg Syndrome |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$MTF1/MTF-1.01</u> | Metal transcription factor 1, MRE |
| <u>V\$SF1F/FTF.01</u> | Alpha (1)-fetoprotein transcription factor (FTF), liver receptor homologue-1 (LHR-1) |
| <u>V\$SMAD/SMAD4.01</u> | Smad4 transcription factor involved in TGF-beta signaling |
| <u>V\$NFKB/NFKAPPAB.01</u> | NF-kappaB |
| <u>V\$EKLF/EKLF.01</u> | Erythroid krueppel like factor (EKLF) |
| <u>V\$CREB/TAXCREB.01</u> | Tax/CREB complex |
| <u>V\$E2FF/E2F.03</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| <u>V\$CP2F/CP2.01</u> | CP2 |
| <u>V\$AHRR/AHRARNT.01</u> | Aryl hydrocarbon receptor / Arnt heterodimers |
| <u>V\$EGRF/EGR2.01</u> | Egr-2/Krox-20 early growth response gene product |
| <u>V\$ZF5F/ZF5.01</u> | Zinc finger / POZ domain transcription factor |
| <u>V\$EBOR/XBP1.01</u> | X-box-binding protein 1 |
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$EGRF/NGFIC.01</u> | Nerve growth factor-induced protein C |
| <u>V\$PCAT/ACAAT.01</u> | Avian C-type LTR CCAAT box |

| Family/matrix | Further Information |
|------------------------------|---|
| <u>V\$PBXC/PBX1 MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$AHRR/AHRARNT.02</u> | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$GREF/GRE.01</u> | Glucocorticoid receptor, C2C2 zinc finger protein binds glucocorticoid dependent to GREs |
| <u>V\$NEUR/NEUROD1.01</u> | DNA binding site for NEUROD1 (BETA-2 / E47 dimer) |
| <u>V\$NRSE/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$NRSE/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$AHRR/AHRARNT.02</u> | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$EGRF/EGR3.01</u> | early growth response gene 3 product |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$AP2F/AP2.01</u> | Activator protein 2 |
| <u>V\$HIF/HIF1.02</u> | Hypoxia inducible factor, bHLH / PAS protein family |
| <u>V\$NRSE/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$ZFIA/ZID.01</u> | zinc finger with interaction domain |
| <u>V\$SMAD/SMAD4.01</u> | Smad4 transcription factor involved in TGF-beta signaling |
| <u>V\$AHRR/AHRARNT.02</u> | Aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$EBOX/MYCMAX.01</u> | c-Myc/Max heterodimer |
| <u>V\$EBOX/USF.03</u> | upstream stimulating factor |
| <u>V\$EGRF/EGR1.01</u> | Egr-1/Krox-24/NGFI-A immediate-early gene product |
| <u>V\$MINI/MUSCLE INI.01</u> | Muscle Initiator Sequence |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$NRSE/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$NF1F/NF1.01</u> | Nuclear factor 1 |
| <u>V\$SF1F/SF1.01</u> | SF1 steroidogenic factor 1 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2A3

After removal of TFBS from hluc+ver2A2 = before removal of TFBS
from hluc+ver2A3 (8 matches)

5

| Family/matrix | Further Information |
|----------------------------|---|
| <u>V\$EGRF/EGR2.01</u> | Egr-2/Krox-20 early growth response gene product |
| <u>V\$HAML/AML3.01</u> | Runt-related transcription factor 2 / CBFA1 (core-binding factor, runt domain, alpha subunit 1) |
| <u>V\$MYOF/MYOGNF1.01</u> | Myogenin / nuclear factor 1 or related factors |
| <u>V\$NF1F/NF1.01</u> | Nuclear factor 1 |
| <u>V\$ETSF/GABP.01</u> | GABP: GA binding protein |
| <u>V\$NFKB/NFKAPPAB.01</u> | NF-kappaB |
| <u>V\$EKLF/EKLF.01</u> | Erythroid krueppel like factor (EKLF) |
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2A6

10

After removal of TFBS from hluc+ver2A5 (2 matches)

| Family/matrix | Further Information |
|------------------------|---|
| <u>V\$HAML/AML3.01</u> | Runt-related transcription factor 2 / CBFA1 (core-binding factor, runt domain, alpha subunit 1) |
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2A6

15

Before removal of TFBS from hluc+ver2A6 (4 matches)

| Family/matrix | Further Information |
|------------------------|--|
| <u>V\$PAX5/PAX5.03</u> | PAX5 paired domain protein |
| <u>V\$LEFF/LEF1.02</u> | TCF/LEF-1, involved in the Wnt signal transduction pathway |
| <u>V\$IRFF/IRF7.01</u> | Interferon regulatory factor 7 (IRF-7) |
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2A7

After removal of TFBS from hluc+ver2A6 = before removal of TFBS from hluc+ver2A7 (1 match)

| Family/matrix | Further Information |
|------------------------|-----------------------------------|
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |

5

TFBS in hluc+ver2A8

After removal of TFBS from hluc+ver2A7 (1 match)

| Family/matrix | Further Information |
|------------------------|-----------------------------------|
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |

10

Table 24

Sequences in Synthetic Luc Genes (version B)15 TFBS in hluc+ver2B1

Before removal of TFBS from hluc+ver2B1 (187 matches)

| Family/matrix | Further Information |
|--------------------------|--|
| <u>V\$HOXF/PTX1.01</u> | Pituitary Homeobox 1 (Ptx1) |
| <u>V\$OCT1/OCT1.04</u> | octamer-binding factor 1 |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$NKXH/NKX25.02</u> | homeo domain factor Nkx-2.5/Csx, tinman homolog low affinity sites |
| <u>V\$BARB/BARBIE.01</u> | barbiturate-inducible element |
| <u>V\$TBPF/TATA.01</u> | cellular and viral TATA box elements |
| <u>V\$GATA/GATA.01</u> | GATA binding site (consensus) |
| <u>V\$AP4R/AP4.01</u> | Activator protein 4 |
| <u>V\$HEN1/HEN1.02</u> | HEN1 |
| <u>V\$SRFF/SRF.01</u> | serum response factor |
| <u>V\$PARF/DBP.01</u> | Albumin D-box binding protein |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$EVI1/EVI1.04</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$GFI1/Gfi1B.01</u> | Growth factor independence 1 zinc finger protein Gfi-1B |
| <u>V\$RBPF/RBPJK.01</u> | Mammalian transcriptional repressor RBP-Jkappa/CBF1 |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |

| Family/matrix | Further Information |
|--------------------------|--|
| V\$AP4R/TAL1ALPHA/E47.01 | Tal-1alpha/E47 heterodimer |
| V\$SRFF/SRF.01 | serum response factor |
| V\$OCTP/OCT1P.01 | octamer-binding factor 1, POU-specific domain |
| V\$BRNF/BRN2.01 | POU factor Brn-2 (N-Oct 3) |
| V\$CREB/E4BP4.01 | E4BP4, bZIP domain, transcriptional repressor |
| V\$VBP/VBP.01 | PAR-type chicken vitellogenin promoter-binding protein |
| V\$EVII/EVII.04 | Ecotropic viral integration site 1 encoded factor |
| V\$CLOX/CDPCR3.01 | cut-like homeodomain protein |
| V\$GFI1/Gfi1B.01 | Growth factor independence 1 zinc finger protein Gfi-1B |
| V\$GATA/LMO2COM.02 | complex of Lmo2 bound to Tal-1, E2A proteins, and GATA-1, half-site 2 |
| V\$SRFF/SRF.01 | serum response factor |
| V\$HOXT/MEIS1 HOXA9.01 | Homeobox protein MEIS1 binding site |
| V\$OCT1/OCT1.03 | octamer-binding factor 1 |
| V\$GFI1/GFI1.01 | Growth factor independence 1 zinc finger protein acts as transcriptional repressor |
| V\$HNF6/HNF6.01 | Liver enriched Cut - Homeodomain transcription factor HNF6 (ONECUT) |
| V\$AML/AML1.01 | runt-factor AML-1 |
| V\$GREF/PRE.01 | Progesterone receptor binding site |
| V\$STAT/STAT5.01 | STAT5: signal transducer and activator of transcription 5 |
| V\$TBPF/TATA.01 | cellular and viral TATA box elements |
| V\$CLOX/CDP.01 | cut-like homeodomain protein |
| V\$FKHD/HFH8.01 | HNF-3/Fkh Homolog-8 |
| V\$FAST/FAST1.01 | FAST-1 SMAD interacting protein |
| V\$GFI1/Gfi1B.01 | Growth factor independence 1 zinc finger protein Gfi-1B |
| V\$CART/CART1.01 | Cart-1 (cartilage homeoprotein 1) |
| V\$HMTB/MTBF.01 | muscle-specific Mt binding site |
| V\$TBPF/TATA.01 | cellular and viral TATA box elements |
| V\$FKHD/XFD2.01 | Xenopus fork head domain factor 2 |
| V\$BRNF/BRN2.01 | POU factor Brn-2 (N-Oct 3) |
| V\$MEF2/AMEF2.01 | myocyte enhancer factor |
| V\$BRNF/BRN2.01 | POU factor Brn-2 (N-Oct 3) |
| V\$BEL1/BEL1.01 | Bel-1 similar region (defined in Lentivirus LTRs) |
| V\$NOLF/OLF1.01 | olfactory neuron-specific factor |

| Family/matrix | Further Information |
|----------------------------|--|
| <u>V\$OCT1/OCT1.06</u> | octamer-binding factor 1 |
| <u>V\$NFKB/NFKAPPAB.02</u> | NF-kappaB |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$HEAT/HSF1.01</u> | heat shock factor 1 |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$PIT1/PIT1.01</u> | Pit1, GHF-1 pituitary specific pou domain transcription factor |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |
| <u>V\$HNF6/HNF6.01</u> | Liver enriched Cut - Homeodomain transcription factor HNF6 (ONECUT) |
| <u>V\$CLOX/CLOX.01</u> | Clox |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$HOXF/PTX1.01</u> | Pituitary Homeobox 1 (Ptx1) |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$FKHD/FREAC4.01</u> | Fork head RElated ACTivator-4 |
| <u>V\$E4FF/E4F.01</u> | GLI-Krueppel-related transcription factor, regulator of adenovirus E4 promoter |
| <u>V\$PDX1/ISL1.01</u> | Pancreatic and intestinal lim-homeodomain factor |
| <u>V\$CART/CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |
| <u>V\$GFI1/GFI1.01</u> | Growth factor independence 1 zinc finger protein acts as transcriptional repressor |
| <u>V\$IRF3/IRF3.01</u> | Interferon regulatory factor 3 (IRF-3) |
| <u>V\$BARB/BARBE.01</u> | barbiturate-inducible element |
| <u>V\$PBXF/PBX1.01</u> | homeo domain factor Pbx-1 |
| <u>V\$EVI1/EVI1.02</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$GATA/GATA2.01</u> | GATA-binding factor 2 |
| <u>V\$BRNF/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$PARF/DBP.01</u> | Albumin D-box binding protein |
| <u>V\$BRNF/BRN3.01</u> | POU transcription factor Brn-3 |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$CREB/TAXCREB.02</u> | Tax/CREB complex |
| <u>V\$GREF/PRE.01</u> | Progesterone receptor binding site |
| <u>V\$RBP/RBPJK.01</u> | Mammalian transcriptional repressor RBP- |

| Family/matrix | Further Information |
|------------------------|--|
| | Jkappa/CBF1 |
| V\$GATA/GATA3.02 | GATA-binding factor 3 |
| V\$STAT/STAT.01 | signal transducers and activators of transcription |
| V\$IKRS/IK2.01 | Ikaros 2, potential regulator of lymphocyte differentiation |
| V\$SRFF/SRF.01 | serum response factor |
| V\$SEF1/SEF1.01 | SEF1 binding site |
| V\$HAML/AML1.01 | runt-factor AML-1 |
| V\$MOKF/MOK2.02 | Ribonucleoprotein associated zinc finger protein MOK-2 (human) |
| V\$FKHD/FREAC2.01 | Fork head RElated ACTivator-2 |
| V\$HMTB/MTBF.01 | muscle-specific Mt binding site |
| V\$GFI1/GFI1.01 | Growth factor independence 1 zinc finger protein acts as transcriptional repressor |
| V\$ECAT/NFY.03 | nuclear factor Y (Y-box binding factor) |
| V\$HOXT/MEIS1 HOXA9.01 | Homeobox protein MEIS1 binding site |
| V\$PCAT/ACAAT.01 | Avian C-type LTR CCAAT box |
| V\$HNF6/HNF6.01 | Liver enriched Cut - Homeodomain transcription factor HNF6 (ONECUT) |
| V\$CLOX/CLOX.01 | Clox |
| V\$GATA/GATA3.02 | GATA-binding factor 3 |
| V\$AREB/AREB6.04 | AREB6 (Atp1a1 regulatory element binding factor 6) |
| V\$GATA/GATA3.02 | GATA-binding factor 3 |
| V\$FKHD/HNF3B.01 | Hepatocyte Nuclear Factor 3beta |
| V\$IRFF/IRF1.01 | interferon regulatory factor 1 |
| V\$NKXH/NKX31.01 | prostate-specific homeodomain protein NKX3.1 |
| V\$PBXF/PBX1.01 | homeo domain factor Pbx-1 |
| V\$ECAT/NFY.03 | nuclear factor Y (Y-box binding factor) |
| V\$PBXC/PBX1 MEIS1.02 | Binding site for a Pbx1/Meis1 heterodimer |
| V\$CLOX/CDP.02 | transcriptional repressor CDP |
| V\$HOXT/MEIS1 HOXA9.01 | Homeobox protein MEIS1 binding site |
| V\$HOXF/HOXA9.01 | Member of the vertebrate HOX - cluster of homeobox factors |
| V\$GATA/GATA.01 | GATA binding site (consensus) |
| V\$NKXH/NKX31.01 | prostate-specific homeodomain protein NKX3.1 |
| V\$GATA/GATA3.02 | GATA-binding factor 3 |
| V\$HOXF/CRX.01 | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |
| V\$CART/CART1.01 | Cart-1 (cartilage homeoprotein 1) |

| Family/matrix | Further information |
|---------------------------------|---|
| <u>V\$OCT1/OCT1.02</u> | octamer-binding factor 1 |
| <u>V\$MAZF/MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |
| <u>V\$CLOX/CDPCR3.01</u> | cut-like homeodomain protein |
| <u>V\$AP1F/VMAF.01</u> | v-Maf |
| <u>V\$AP4R/TAL1ALPHA/E47.01</u> | Tal-1alpha/E47 heterodimer |
| <u>V\$PAX8/PAX8.01</u> | PAX 2/5/8 binding site |
| <u>V\$BRAC/BRACH.01</u> | Brachyury |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$RREB/RREB1.01</u> | Ras-responsive element binding protein 1 |
| <u>V\$MZF1/MZF1.01</u> | MZF1 |
| <u>V\$MOKF/MOK2.02</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (human) |
| <u>V\$HOXF/PTX1.01</u> | Pituitary Homeobox 1 (Ptx1) |
| <u>V\$LTUP/TAACC.01</u> | Lentiviral TATA upstream element |
| <u>V\$AP4R/TH1E47.01</u> | Thing1/E47 heterodimer, TH1 bHLH member specific expression in a variety of embryonic tissues |
| <u>V\$XSEC/STAF.01</u> | Se-Cys tRNA gene transcription activating factor |
| <u>V\$IKRS/IK3.01</u> | Ikaros 3, potential regulator of lymphocyte differentiation |
| <u>V\$AP1F/AP1.01</u> | AP1 binding site |
| <u>V\$MAZF/MAZ.01</u> | Myc associated zinc finger protein (MAZ) |
| <u>V\$MZF1/MZF1.01</u> | MZF1 |
| <u>V\$CLOX/CDPCR3.01</u> | cut-like homeodomain protein |
| <u>V\$P53F/P53.01</u> | tumor suppressor p53 |
| <u>V\$SMAD/SMAD3.01</u> | Smad3 transcription factor involved in TGF-beta signaling |
| <u>V\$HMTB/MTBF.01</u> | muscle-specific Mt binding site |
| <u>V\$OCT1/OCT1.03</u> | octamer-binding factor 1 |
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |
| <u>V\$PIT1/PIT1.01</u> | Pit1, GHF-1 pituitary specific pou domain transcription factor |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$HOXF/HOX1-3.01</u> | Hox-1.3, vertebrate homeobox protein |

| Family/matrix | Further Information |
|-------------------------------|--|
| <u>V\$PBXF/PBX1.01</u> | homeo domain factor Pbx-1 |
| <u>V\$ECAT/NFY.03</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$PBXC/PBX1 MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$CLOX/CDP.02</u> | transcriptional repressor CDP |
| <u>V\$HOXT/MEIS1 HOXA9.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$HOXF/HOXA9.01</u> | Member of the vertebrate HOX - cluster of homeobox factors |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$PCAT/ACAAT.01</u> | Avian C-type LTR CCAAT box |
| <u>V\$XSEC/STAF.01</u> | Se-Cys tRNA gene transcription activating factor |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$CLOX/CDP.01</u> | cut-like homeodomain protein |
| <u>V\$FAST/FAST1.01</u> | FAST-1 SMAD interacting protein |
| <u>V\$ECAT/NFY.01</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$MEF2/MMEF2.01</u> | myocyte enhancer factor |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |
| <u>V\$FAST/FAST1.01</u> | FAST-1 SMAD interacting protein |
| <u>V\$LTUP/TAACC.01</u> | Lentiviral TATA upstream element |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$BRNF/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$HEN1/HEN1.01</u> | HEN 1 |
| <u>V\$BEL1/BEL1.01</u> | Bel-1 similar region (defined in Lentivirus LTRs) |
| <u>V\$HOXF/PTX1.01</u> | Pituitary Homeobox 1 (Ptx1) |
| <u>V\$BRNF/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$NFKB/NFKAPPAB.01</u> | NF-kappaB |
| <u>V\$HAML/AML1.01</u> | runt-factor AML-1 |
| <u>V\$ZFIA/ZID.01</u> | zinc finger with interaction domain |
| <u>V\$XSEC/STAF.02</u> | Se-Cys tRNA gene transcription activating factor |
| <u>V\$IKRS/IK1.01</u> | Ikaros 1, potential regulator of lymphocyte differentiation |
| <u>V\$FAST/FAST1.01</u> | FAST-1 SMAD interacting protein |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |

| Family/matrix | Further Information |
|-------------------------|--|
| <u>V\$BEL1/BEL1.01</u> | Bel-1 similar region (defined in Lentivirus LTRs) |
| <u>V\$EGRF/WT1.01</u> | Wilms Tumor Suppressor |
| <u>V\$MAZF/MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$ZBPF/ZBP89.01</u> | Zinc finger transcription factor ZBP-89 |
| <u>V\$SPIF/GC.01</u> | GC box elements |
| <u>V\$RREB/RREB1.01</u> | Ras-responsive element binding protein 1 |
| <u>V\$MOKF/MOK2.01</u> | Ribonucleoprotein associated zinc finger protein MOK-2 (mouse) |
| <u>V\$MEIS/MEIS1.01</u> | Binding site for monomeric Meis 1 homeodomain protein |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |
| <u>V\$MAZF/MAZR.01</u> | MYC-associated zinc finger protein related transcription factor |
| <u>V\$MZF1/MZF1.01</u> | MZF1 |
| <u>V\$PDX1/PDX1.01</u> | Pdx1 (IDX1/IPF1) pancreatic and intestinal homeodomain TF |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2B3

- After removal of TFBS from hluc+ver2B2 = before removal of TFBS
 5 from hluc+ver2B3 (35 matches)

| Family/matrix | Further Information |
|----------------------------|--|
| <u>V\$OCT1/OCT1.04</u> | octamer-binding factor 1 |
| <u>V\$BARB/BARBE.01</u> | barbiturate-inducible element |
| <u>V\$NFKB/NFKAPPAB.02</u> | NF-kappaB |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$PIT1/PIT1.01</u> | Pit1, GHF-1 pituitary specific pou domain transcription factor |
| <u>V\$HOXF/PTX1.01</u> | Pituitary Homeobox 1 (Ptx1) |
| <u>V\$FKHD/FREAC4.01</u> | Fork head RElated ACTivator-4 |

| Family/matrix | Further information |
|---------------------------------|--|
| <u>V\$E4FF/E4F.01</u> | GLI-Krueppel-related transcription factor, regulator of adenovirus E4 promoter |
| <u>V\$EVI1/EVI1.02</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$GATA/GATA2.01</u> | GATA-binding factor 2 |
| <u>V\$GREF/PRE.01</u> | Progesterone receptor binding site |
| <u>V\$RBPF/RBPJK.01</u> | Mammalian transcriptional repressor RBP-Jkappa/CBF1 |
| <u>V\$STAT/STAT.01</u> | signal transducers and activators of transcription |
| <u>V\$IKRS/IK2.01</u> | Ikars 2, potential regulator of lymphocyte differentiation |
| <u>V\$FKHD/FREAC2.01</u> | Fork head RElated ACTivator-2 |
| <u>V\$SRFF/SRF.01</u> | serum response factor |
| <u>V\$GREF/PRE.01</u> | Progesterone receptor binding site |
| <u>V\$CLOX/CDPCR3.01</u> | cut-like homeodomain protein |
| <u>V\$AP4R/TAL1ALPHA/E47.01</u> | Tal-1alpha/E47 heterodimer |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$FKHD/XFD3.01</u> | Xenopus fork head domain factor 3 |
| <u>V\$PBXF/PBX1.01</u> | homeo domain factor Pbx-1 |
| <u>V\$ECAT/NFY.03</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$PBXC/PBX1 MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$CLOX/CDP.02</u> | transcriptional repressor CDP |
| <u>V\$HOXT/MEIS1 HOXA9.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$HOXF/HOXA9.01</u> | Member of the vertebrate HOX - cluster of homeobox factors |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$MINI/MUSCLE INI.01</u> | Muscle Initiator Sequence |
| <u>V\$CLOX/CDP.01</u> | cut-like homeodomain protein |
| <u>V\$BRNF/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$NFKB/NFKAPPAB.01</u> | NF-kappaB |
| <u>V\$ZFIA/ZID.01</u> | zinc finger with interaction domain |
| <u>V\$BCL6/BCL6.02</u> | POZ/zinc finger protein, transcriptional repressor, translocations observed in diffuse large cell lymphoma |
| <u>V\$HOXF/CRX.01</u> | Cone-rod homeobox-containing transcription factor / otx-like homeobox gene |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2B6

After removal of TFBS from hluc+ver2B5 (2 matches)

| Family/matrix | Further Information |
|-----------------|-----------------------------------|
| V\$HOXF/PTX1.01 | Pituitary Homeobox 1 (Ptx1) |
| V\$FKHD/XFD3.01 | Xenopus fork head domain factor 3 |

**matches are listed in order of occurrence in the corresponding sequence

5

TFBS in hluc+ver2B6

Before removal of TFBS from hluc+ver2B6 (6 matches)

| Family/matrix | Further Information |
|--------------------|--|
| V\$PAX6/PAX4_PD.01 | PAX4 paired domain binding site |
| V\$HOXF/PTX1.01 | Pituitary Homeobox 1 (Ptx1) |
| V\$FKHD/XFD3.01 | Xenopus fork head domain factor 3 |
| V\$PAX6/PAX6.02 | PAX6 paired domain and homeodomain are required for binding to this site |
| V\$PAX5/PAX5.03 | PAX5 paired domain protein |
| V\$IRFF/IRF3.01 | Interferon regulatory factor 3 (IRF-3) |

**matches are listed in order of occurrence in the corresponding sequence

10

TFBS in hluc+ver2B7

After removal of TFBS from hluc+ver2B6 = before removal of TFBS from hluc+ver2B7 (2 matches)

| Family/matrix | Further Information |
|-----------------|-----------------------------------|
| V\$HOXF/PTX1.01 | Pituitary Homeobox 1 (Ptx1) |
| V\$FKHD/XFD3.01 | Xenopus fork head domain factor 3 |

15 **matches are listed in order of occurrence in the corresponding sequence

TFBS in hluc+ver2B8

After removal of TFBS from hluc+ver2B7 = before removal of TFBS from hluc+ver2B8 (1 match)

20

| Family/matrix | Further Information |
|-----------------|-----------------------------------|
| V\$FKHD/XFD3.01 | Xenopus fork head domain factor 3 |

TFBS in hluc+ver2B9

After removal of TFBS from hluc+ver2B8 = before removal of TFBS from hluc+ver2B9 (1 match)

| Family/matrix | Further Information |
|-----------------|-----------------------------------|
| V\$FKHD/XFD3.01 | Xenopus fork head domain factor 3 |

TFBS in hluc+ver2B10

- 5 After removal of TFBS from hluc+ver2B9 (1 match)

| Family/matrix | Further Information |
|-----------------|-----------------------------------|
| V\$FKHD/XFD3.01 | Xenopus fork head domain factor 3 |

Example 8Summary of Design for pGL4 Sequences

- 10 Figure 2 depicts the design scheme for the pGL4 vector. A portion of the vector backbone in pGL3 which includes an *bla* gene and a sequence between *bla* and a multiple cloning region, but not a second open reading frame, was modified to yield pGL4. pGL4 includes an ampicillin resistance gene between a *NotI* and a *SpeI* site, the sequence of which was modified to remove regulatory
- 15 sequences but not to optimize codons for mammalian expression (*bla*-1-*bla*-5), and a *SpeI*-*NcoI* fragment that includes a multiple cloning region and a translation trap. The translation trap includes about 60 nucleotides having at least two stop codons in each reading frame. The *SpeI*-*NcoI* fragment from a parent vector, pGL4-basics-5F2G-2, was modified to decrease undesired
- 20 regulatory sequences (MCS-1 to MCS-4; SEQ ID Nos. 76-79). One of the resulting sequences, MCS-4, was combined with a modified ampicillin resistance gene, *bla*-5 (SEQ ID NO:84), to yield pGL4B-4NN (SEQ ID NO:95). pGL4B-4NN was further modified (pGL4-NN1-3; SEQ ID Nos. 96-98). To determine if additional polyA sequences in the *SpeI*-*NcoI* fragment further
- 25 reduced expression from the vector backbone, various polyA sequences were inserted therein. For instance, pGL4NN-Blue Heron included a *c-mos* polyA sequence in the *SpeI*-*NcoI* fragment. However, removal of regulatory sequences in polyA sequences may alter the secondary structure and thus the function of those sequences.

In one vector, the *SpeI-NcoI* fragment from pGL3 (*SpeI-NcoI* start ver 2; SEQ ID NO:48) was modified to remove one transcription factor binding site and one restriction enzyme recognition site, and alter the multiple cloning region, yielding *SpeI-NcoI* ver2 (SEQ ID NO:49).

5

TF binding sites and search parameters

Each TF binding site ("matrix") belongs to a matrix family that groups functionally similar matrices together, eliminating redundant matches by MatInspector professional (the search program). Searches were limited to vertebrate TF binding sites. Searches were performed by matrix family, i.e., the results show only the best match from a family for each site. MatInspector default parameters were used for the core and matrix similarity values (core similarity = 0.75, matrix similarity = optimized), except for sequence MCS-1 (core similarity = 1.00, matrix similarity = optimized).

15

Table 25

Description of Designed Sequences

pGL4 sequences

| Sequence | Description | Matrix Library |
|----------|---|------------------|
| | <i>SpeI-NcoI</i> fragment with MCS, translation trap | |
| MCS-1 | <i>SpeI-NcoI</i> from pGL4-basics-5F2G-2 | Ver 2.2 Sep 2001 |
| MCS-2 | First removal of undesired sequence matches | Ver 2.2 Sep 2001 |
| MCS-3 | Second removal of undesired sequence matches | Ver 2.2 Sep 2001 |
| MCS-4 | Third removal of undesired sequence matches | Ver 2.3 Feb 2001 |
| | <i>NotI-SpeI</i> fragment with <i>bla</i> gene | |
| Bla | Beta-lactamase gene from pGL3 vectors | |
| bla-1* | <i>SacII</i> (RE) added, <i>BsmAI</i> (RE) site removed (*) | Ver 2.2 Sep 2001 |
| bla-2* | First removal of undesired sequence matches | Ver 2.3 Feb 2001 |
| bla-3* | Second removal of undesired sequence matches | Ver 2.3 Feb 2001 |
| bla-4* | Third removal of undesired sequence matches | Ver 2.3 Feb 2001 |

| Sequence | Description | Matrix Library |
|-----------------------------|---|------------------|
| bla-5* | Fourth removal of undesired sequence matches | Ver 2.3 Feb 2001 |
| | <i>NotI-NcoI fragment with bla, translation trap, MCS</i> | |
| pGL4B-4NN | Combination of bla-5 and MCS-4 sections | Ver 2.4 May 2002 |
| pGL4B-4NN1 | First removal of undesired sequence matches | Ver 2.4 May 2002 |
| pGL4B-4NN2 | Second removal of undesired sequence matches | Ver 2.4 May 2002 |
| pGL4B-4NN3 | Third version after removal of CEBP (TF) site | Ver 2.4 May 2002 |
| | <i>SpeI-NcoI fragment with translation trap, polyA, MCS</i> | |
| <i>SpeI-NcoI-Ver2-start</i> | Existing MCS replaced with new MCS | Ver 4.0 Nov 2003 |
| <i>SpeI-NcoI-Ver2</i> | First removal of undesired sequence matches | Ver 4.0 Nov 2003 |

(*)Bla codon usage was not optimized for expression in mammalian cells. Low usage *E. coli* codons were avoided when changes were introduced to remove undesired sequence elements.

5

Table 26

Sequences in Synthetic *SpeI-NcoI* fragment of pGL4

TFBS in MCS-1

Before removal of TFBS from MCS-1 (14 matches)

| Name of family/matrix ** | Further Information |
|--------------------------|--|
| <u>V\$PAX3/PAX3.01</u> | Pax-3 paired domain protein, expressed in embryogenesis, mutations correlate to Waardenburg Syndrome |
| <u>V\$GATA/GATA.01</u> | GATA binding site (consensus) |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein |

| | |
|-------------------------|---|
| | NKX3.1 |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$BRN2/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$ZFIA/ZID.01</u> | zinc finger with interaction domain |
| <u>V\$CP2F/CP2.01</u> | CP2 |
| <u>V\$BRAC/BRACH.01</u> | Brachyury |
| <u>V\$PAX6/PAX6.01</u> | Pax-6 paired domain protein |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$TEAF/TEF1.01</u> | TEF-1 related muscle factor |
| <u>V\$ETSF/ELK1.02</u> | Elk-1 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in MCS-2

After removal of TFBS from MCS-1 = before removal of TFBS from MCS-2

5 (12 matches)

| Name of family/matrix ** | Further Information |
|--------------------------|---------------------|
|--------------------------|---------------------|

| | |
|-------------------------|---|
| <u>V\$GATA/GATA.01</u> | GATA binding site (consensus) |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$CART/CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$BRN2/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$PAX6/PAX6.01</u> | Pax-6 paired domain protein |
| <u>V\$PAX8/PAX8.01</u> | PAX 2/5/8 binding site |
| <u>V\$PAX1/PAX1.01</u> | Pax1 paired domain protein, expressed in the developing vertebral column of mouse embryos |

****matches are listed in order of occurrence in the corresponding sequence**

TFBS in MCS-3

After removal of TFBS from MCS-2 = before removal of TFBS from MCS-4

5 (0 matches)

TFBS in MCS-4

After removal of TFBS from MCS-3 (0 matches)

Table 27

Sequences in Synthetic *NotI-SpeI* Fragment of pGL4

TFBS in bla-1

Before removal of TFBS from bla-1 (94 matches)

| Name of family/matrix | Further Information |
|--------------------------|---|
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$HOXF/HOX1-3.01</u> | Hox-1.3, vertebrate homeobox protein |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$SETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$SETSF/ELK1.02</u> | Elk-1 |
| <u>V\$GKLF/GKLF.01</u> | gut-enriched Krueppel-like factor |
| <u>V\$E2FF/E2F.02</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| <u>V\$SETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$AP1F/VMAF.01</u> | v-Maf |
| <u>V\$XBBF/RFX1.01</u> | X-box binding protein RFX1 |
| <u>V\$AREB/AREB6.04</u> | AREB6 (Atp1a1 regulatory element binding factor 6) |
| <u>V\$CMYB/CMYB.01</u> | c-Myb, important in hematopoiesis, cellular equivalent to avian myoblastosis virus oncogene v-myb |
| <u>V\$VMYB/VMYB.02</u> | v-Myb |
| <u>V\$EBOX/NMYC.01</u> | N-Myc |
| <u>V\$VBPF/VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |
| <u>V\$CMYB/CMYB.01</u> | c-Myb, important in hematopoiesis, cellular equivalent to avian |

| Name of family/matrix | Further Information |
|-------------------------|---|
| | myoblastosis virus oncogene v-myb |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |
| <u>V\$PAX8/PAX8.01</u> | PAX 2/5/8 binding site |
| <u>V\$HNF4/HNF4.02</u> | Hepatic nuclear factor 4 |
| <u>V\$E2FF/E2F.01</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| <u>V\$NFAT/NFAT.01</u> | Nuclear factor of activated T-cells |
| <u>V\$ECAT/NFY.02</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |
| <u>V\$MYT1/MYT1.02</u> | MyT1 zinc finger transcription factor involved in primary neurogenesis |
| <u>V\$GATA/GATA3.01</u> | GATA-binding factor 3 |
| <u>V\$CREB/CREB.02</u> | cAMP-responsive element binding protein |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$IRFF/ISRE.01</u> | interferon-stimulated response element |
| <u>V\$NRSF/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$TCFF/TCF11.01</u> | TCF11/KCR-F1/Nrf1 homodimers |
| <u>V\$STAT/STAT.01</u> | signal transducers and activators of transcription |
| <u>V\$ECAT/NFY.03</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$OCT1/OCT1.05</u> | octamer-binding factor 1 |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$NKXH/NKX25.02</u> | homeo domain factor Nkx-2.5/Csx, |

| Name of family/matrix | Further Information |
|------------------------------|--|
| | tinman homolog low affinity sites |
| <u>V\$PIT1/PIT1.01</u> | Pit1, GHF-1 pituitary specific pou domain transcription factor |
| <u>V\$CLOX/CDPCR3.01</u> | cut-like homeodomain protein |
| <u>V\$GREF/ARE.01</u> | Androgene receptor binding site |
| <u>V\$GATA/GATA1.04</u> | GATA-binding factor 1 |
| <u>V\$E2TF/E2.02</u> | papilloma virus regulator E2 |
| <u>V\$RPOA/POLYA.01</u> | Mammalian C-type LTR Poly A signal |
| <u>V\$VMYB/VMYB.02</u> | v-Myb |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$VBPF/VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |
| <u>V\$CREB/HLF.01</u> | hepatic leukemia factor |
| <u>V\$SF1F/SF1.01</u> | SF1 steroidogenic factor 1 |
| <u>V\$XBBF/MIF1.01</u> | MIBP-1 / RFX1 complex |
| <u>V\$IKRS/TK2.01</u> | Ikaros 2, potential regulator of lymphocyte differentiation |
| <u>V\$MINI/MUSCLE INI.02</u> | Muscle Initiator Sequence |
| <u>V\$PCAT/CLTR CAAT.01</u> | Mammalian C-type LTR CCAAT box |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$RPAD/PADS.01</u> | Mammalian C-type LTR Poly A downstream element |
| <u>V\$XBBF/RFX1.02</u> | X-box binding protein RFX1 |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$CREB/HLF.01</u> | hepatic leukemia factor |
| <u>V\$HNF1/HNF1.01</u> | hepatic nuclear factor 1 |

| <u>Name of family/matrix</u> | <u>Further Information</u> |
|------------------------------|---|
| <u>V\$VMYB/VMYB.01</u> | v-Myb |
| <u>V\$NKXH/NKX31.01</u> | prostate-specific homeodomain protein NKX3.1 |
| <u>V\$XBBF/RFX1.01</u> | X-box binding protein RFX1 |
| <u>V\$STAT/STAT.01</u> | signal transducers and activators of transcription |
| <u>V\$HNF1/HNF1.01</u> | hepatic nuclear factor 1 |
| <u>V\$HMYO/S8.01</u> | S8 |
| <u>V\$SORY/SOX5.01</u> | Sox-5 |
| <u>V\$RBIT/BRIGHT.01</u> | Bright, B cell regulator of IgH transcription |
| <u>V\$NKXH/NKX25.02</u> | homeo domain factor Nkx-2.5/Csx, tinman homolog low affinity sites |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$BARB/BARBIE.01</u> | barbiturate-inducible element |
| <u>V\$MTF1/MTF-1.01</u> | Metal transcription factor 1, MRE |
| <u>V\$NFKB/CREL.01</u> | c-Rel |
| <u>V\$ETSF/ELK1.02</u> | Elk-1 |
| <u>V\$CLOX/CDP.01</u> | cut-like homeodomain protein |
| <u>V\$RPOA/LPOLYA.01</u> | Lentiviral Poly A signal |
| <u>V\$GATA/GATA1.03</u> | GATA-binding factor 1 |
| <u>V\$ZFIA/ZID.01</u> | zinc finger with interaction domain |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$PAX1/PAX1.01</u> | Pax1 paired domain protein, expressed in the developing vertebral column of |

| <u>Name of family/matrix</u> | <u>Further Information</u> |
|-------------------------------|--|
| | mouse embryos |
| <u>V\$GATA/LMO2COM.02</u> | complex of Lmo2 bound to Tal-1, E2A proteins, and GATA-1, half-site 2 |
| <u>V\$NRSF/NRSF.01</u> | neuron-restrictive silencer factor |
| <u>V\$AP4R/TAL1BETAE47.01</u> | Tal-1beta/E47 heterodimer |
| <u>V\$GATA/LMO2COM.02</u> | complex of Lmo2 bound to Tal-1, E2A proteins, and GATA-1, half-site 2 |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$XBBF/RFX1.01</u> | X-box binding protein RFX1 |
| <u>V\$AHRR/AHRARNT.02</u> | aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$PAX5/PAX9.01</u> | zebrafish PAX9 binding sites |
| <u>V\$CLOX/CDP.02</u> | transcriptional repressor CDP |
| <u>V\$GATA/GATA1.01</u> | GATA-binding factor 1 |
| <u>V\$AP1F/TCF11MAFG.01</u> | TCF11/MafG heterodimers, binding to subclass of AP1 sites |
| <u>V\$BRN2/BRN2.01</u> | POU factor Brn-2 (N-Oct 3) |
| <u>V\$NKXH/NKX25.02</u> | homeo domain factor Nkx-2.5/Csx, tinman homolog low affinity sites |
| <u>V\$ECAT/NFY.02</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$FKHD/FREAC4.01</u> | Fork head RElated ACTivator-4 |
| <u>V\$NFAT/NFAT.01</u> | Nuclear factor of activated T-cells |
| <u>V\$IRFF/IRF1.01</u> | interferon regulatory factor 1 |
| <u>V\$E2FF/E2F.02</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in bla-2

After removal of TFBS from bla-1 = before removal of TFBS from bla-2
 = (51 matches)

| <u>Name of family/matrix</u> | <u>Further Information</u> |
|------------------------------|---|
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$ETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$OCTP/OCT1P.01</u> | octamer-binding factor 1, POU-specific domain |
| <u>V\$ETSF/ELK1.02</u> | Elk-1 |
| <u>V\$EBOX/NMYC.01</u> | N-Myc |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |
| <u>V\$PAX8/PAX8.01</u> | PAX 2/5/8 binding site |
| <u>V\$HNF4/HNF4.02</u> | Hepatic nuclear factor 4 |
| <u>V\$E2FF/E2F.01</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| <u>V\$NFAT/NFAT.01</u> | Nuclear factor of activated T-cells |
| <u>V\$ECAT/NFY.02</u> | nuclear factor Y (Y-box binding factor) |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |
| <u>V\$MYT1/MYT1.02</u> | MyT1 zinc finger transcription factor involved in primary neurogenesis |
| <u>V\$GATA/GATA3.01</u> | GATA-binding factor 3 |
| <u>V\$CREB/CREB.02</u> | cAMP-responsive element binding protein |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$NRSE/NRSE.01</u> | neural-restrictive-silencer-element |
| <u>V\$OCT1/OCT1.05</u> | octamer-binding factor 1 |
| <u>V\$CLOX/CDPCR3.01</u> | cut-like homeodomain protein |

| Name of family/matrix | Further Information |
|-------------------------|--|
| <u>V\$GREF/ARE.01</u> | Androgene receptor binding site |
| <u>V\$GATA/GATA1.04</u> | GATA-binding factor 1 |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$CREB/HLF.01</u> | hepatic leukemia factor |
| <u>V\$VBPF/VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |
| <u>V\$XBBF/MIF1.01</u> | MIBP-1 / RFX1 complex |
| <u>V\$IKRS/IK2.01</u> | Ikars 2, potential regulator of lymphocyte differentiation |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$XBBF/RFX1.02</u> | X-box binding protein RFX1 |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$CREB/HLF.01</u> | hepatic leukemia factor |
| <u>V\$XBBF/RFX1.02</u> | X-box binding protein RFX1 |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$BARB/BARBE.01</u> | barbiturate-inducible element |
| <u>V\$MTF1/MTF-1.01</u> | Metal transcription factor 1, MRE |
| <u>V\$NFKB/CREL.01</u> | c-Rel |
| <u>V\$ETSF/ELK1.02</u> | Elk-1 |
| <u>V\$TBPF/TATA.01</u> | cellular and viral TATA box elements |
| <u>V\$MEIS/MEIS1.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$HOXF/HOXA9.01</u> | Member of the vertebrate HOX - cluster of homeobox factors |
| <u>V\$GATA/GATA1.03</u> | GATA-binding factor 1 |
| <u>V\$MEIS/MEIS1.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$NOLF/OLF1.01</u> | olfactory neuron-specific factor |

| Name of family/matrix | Further Information |
|--------------------------------|---|
| <u>V\$AP4R/TAL1BETA.E47.01</u> | Tal-1beta/E47 heterodimer |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$XBBF/RFX1.01</u> | X-box binding protein RFX1 |
| <u>V\$AHRR/AHRARNT.02</u> | aryl hydrocarbon / Arnt heterodimers, fixed core |
| <u>V\$PAX5/PAX9.01</u> | zebrafish PAX9 binding sites |
| <u>V\$CLOX/CDP.02</u> | transcriptional repressor CDP |
| <u>V\$GATA/GATA1.01</u> | GATA-binding factor 1 |
| <u>V\$IRFF/IRF1.01</u> | interferon regulatory factor 1 |
| <u>V\$E2FF/E2F.02</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in bla-3

After removal of TFBS from bla-2 = before removal of TFBS from bla-3

5 = (16 matches)

| Name of family/matrix | Further Information |
|-------------------------|---|
| <u>V\$SETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$E2FF/E2F.02</u> | E2F, involved in cell cycle regulation, interacts with Rb p107 protein |
| <u>V\$NFAT/NFAT.01</u> | Nuclear factor of activated T-cells |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |
| <u>V\$MYT1/MYT1.02</u> | MyT1 zinc finger transcription factor involved in primary neurogenesis |

| Name of family/matrix | Further Information |
|-----------------------------|---|
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$SORY/SOX5.01</u> | Sox-5 |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$CREB/HLF.01</u> | hepatic leukemia factor |
| <u>V\$BPF/VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$XBBF/RFX1.02</u> | X-box binding protein RFX1 |
| <u>V\$CREB/HLF.01</u> | hepatic leukemia factor |
| <u>V\$GATA/GATA1.0</u> 3 | GATA-binding factor 1 |
| <u>V\$MEIS/MEIS1.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$NOLF/OLF1.01</u> | olfactory neuron-specific factor |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in bla-4

After removal of TFBS from bla-3 = before removal of TFBS from bla-4

5 = (14 matches)

| Name of family/matrix** | Further Information |
|-------------------------|---------------------|
|-------------------------|---------------------|

| Name of family/matrix | Further Information |
|---------------------------|---|
| <u>V\$ETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$NFAT/NFAT.01</u> | Nuclear factor of activated T-cells |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$GATA/GATA3.01</u> | GATA-binding factor 3 |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$EBOX/USF.02</u> | upstream stimulating factor |
| <u>V\$PAX5/PAX5.01</u> | B-cell-specific activating protein |
| <u>V\$XBBF/RFX1.02</u> | X-box binding protein RFX1 |
| <u>V\$GATA/GATA1.03</u> | GATA-binding factor 1 |
| <u>V\$MEIS/MEIS1.01</u> | Homeobox protein MEIS1 binding site |
| <u>V\$ZFIA/ZID.01</u> | zinc finger with interaction domain |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$PAX1/PAX1.01</u> | Pax1 paired domain protein, expressed in the developing vertebral column of mouse embryos |
| <u>V\$GATA/LMO2COM.02</u> | complex of Lmo2 bound to Tal-1, E2A |

| Name of family/matrix | Further Information |
|-----------------------|-----------------------------------|
| | proteins, and GATA-1, half-site 2 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in bla-5

After removal of TFBS from bla-4 (5 matches)

| Name of family/matrix | Further Information |
|-------------------------|---|
| <u>V\$ETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$GATA/GATA3.01</u> | GATA-binding factor 3 |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$EBOX/USF.02</u> | upstream stimulating factor |

5 **matches are listed in order of occurrence in the corresponding sequence

Table 28

Sequences in Synthetic *NotI*-*NcoI* Fragment of pGL4

TFBS in pGL4B-4NN

10 Before removal of TFBS from pGL4B-4NN = (11 matches)

| Name of family/matrix** | Further Information |
|-------------------------|---------------------------------|
| <u>V\$SMAD/FAST1.01</u> | FAST-1 SMAD interacting protein |
| <u>V\$SMAD/FAST1.01</u> | FAST-1 SMAD interacting protein |

| | |
|-------------------------|---|
| <u>V\$ETSF/FLI.01</u> | ETS family member FLI |
| <u>V\$RBPF/RBPJK.01</u> | Mammalian transcriptional repressor RBP-Jkappa/CBF1 |
| <u>V\$ETSF/FLI.01</u> | ETS family member FLI |
| <u>V\$EBOX/USF.02</u> | upstream stimulating factor |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$GATA/GATA3.01</u> | GATA-binding factor 3 |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$ETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in pGL4B-4NN1

After removal of TFBS from pGL4B-4NN = before removal of TFBS

5 from pGL4B-4NN1 (7 matches)

| Name of family/matrix | Further Information |
|------------------------------|---|
| <u>V\$ETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |

| | |
|-------------------------|---------------------------------|
| <u>V\$EBOX/USF.02</u> | upstream stimulating factor |
| <u>V\$ETSF/FLI.01</u> | ETS family member FLI |
| <u>V\$SMAD/FAST1.01</u> | FAST-1 SMAD interacting protein |
| <u>V\$SMAD/FAST1.01</u> | FAST-1 SMAD interacting protein |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in pGL4B-4NN2

After removal of TFBS from pGL4B-4NN1 = before removal of TFBS

5 from pGL4B-4NN2 (4 matches)

| Name of family/matrix | Further Information |
|----------------------------------|---|
| <u>V\$ETSF/NRF2.01</u> | nuclear respiratory factor 2 |
| <u>V\$WHZF/WHN.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$CEBP/CEBPB.01</u> | CCAAT/enhancer binding protein beta |
| <u>V\$EBOX/USF.02</u> | upstream stimulating factor |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in pGL4B-4NN3

After removal of TFBS from pGL4B-4NN2 (3 matches)

| Name of family/matrix ** | Further Information |
|---|-----------------------------|
| <u>V\$EBOX/USF.</u> | upstream stimulating factor |

| | |
|-----------------------------------|---|
| <u>02</u> | |
| <u>V\$WHZF/WH</u> <u>N.01</u> | winged helix protein, involved in hair keratinization and thymus epithelium differentiation |
| <u>V\$ETSF/NRF2</u> <u>.01</u> | nuclear respiratory factor 2 |

**matches are listed in order of occurrence in the corresponding sequence

Table 29

Sequences in Synthetic *SpeI-NcoI* section of pGL4

5 TFBS in *SpeI-NcoI*-Ver2-start

Before removal of TFBS from *SpeI-NcoI*-Ver2-start (34 matches)

| <u>Family/matrix</u> | <u>Further Information</u> |
|-------------------------|--|
| <u>V\$PAX8/PAX8.01</u> | PAX 2/5/8 binding site |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$NKXH/NKX31.01</u> | Prostate-specific homeodomain protein NKX3.1 |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$NKXH/NKX31.01</u> | Prostate-specific homeodomain protein NKX3.1 |
| <u>V\$CART/CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |
| <u>V\$NKXH/NKX25.02</u> | Homeo domain factor Nkx-2.5/Csx, tinman homolog low affinity sites |
| <u>V\$ETSF/ELK1.01</u> | Elk-1 |

| Family/matrix | Further Information |
|---------------------------------|---|
| <u>V\$CDXF/CDX2.01</u> | Cdx-2 mammalian caudal related intestinal transcr. factor |
| <u>V\$BRNF/BRN3.01</u> | POU transcription factor Brn-3 |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |
| <u>V\$FKHD/FREAC3.01</u> | Fork head related activator-3 (FOXC1) |
| <u>V\$OCT1/OCT1.02</u> | Octamer-binding factor 1 |
| <u>V\$CART/CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |
| <u>V\$PDX1/PDX1.01</u> | Pdx1 (IDX1/IPF1) pancreatic and intestinal homeodomain TF |
| <u>V\$PARF/DBP.01</u> | Albumin D-box binding protein |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |
| <u>V\$VBPF/VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |
| <u>V\$AP4R/TAL1ALPHA E47.01</u> | Tal-1alpha/E47 heterodimer |
| <u>V\$RP58/RP58.01</u> | Zinc finger protein RP58 (ZNF238), associated preferentially with heterochromatin |
| <u>V\$COMP/COMP1.01</u> | COMP1, cooperates with myogenic proteins in multicomponent complex |
| <u>V\$CLOX/CLOX.01</u> | Clox |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$PBXC/PBX1 MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$PBXF/PBX1.01</u> | Homeo domain factor Pbx-1 |
| <u>V\$IRFF/IRF1.01</u> | Interferon regulatory factor 1 |
| <u>V\$TEAF/TEF1.01</u> | TEF-1 related muscle factor |

| Family/matrix | Further information |
|-------------------------|---|
| <u>V\$EBOX/ATF6.01</u> | Member of b-zip family, induced by ER damage/stress, binds to the ERSE in association with NF-Y |
| <u>V\$NKXH/NKX32.01</u> | Homeodomain protein NKX3.2 (BAPX1, NKX3B, Bagpipe homolog) |
| <u>V\$E2TF/E2.02</u> | Papilloma virus regulator E2 |
| <u>V\$EV11/EV11.05</u> | Ecotropic viral integration site 1 encoded factor |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |

**matches are listed in order of occurrence in the corresponding sequence

TFBS in *SpeI-NcoI-Ver2*

After removal of TFBS from *SpeI-NcoI-Ver2*-start (28 matches)

| Family/matrix | Further information |
|-------------------------|---|
| <u>V\$PAX8/PAX8.01</u> | PAX 2/5/8 binding site |
| <u>V\$GATA/GATA1.02</u> | GATA-binding factor 1 |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$NKXH/NKX31.01</u> | Prostate-specific homeodomain protein NKX3.1 |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$CREB/E4BP4.01</u> | E4BP4, bZIP domain, transcriptional repressor |
| <u>V\$NKXH/NKX31.01</u> | Prostate-specific homeodomain protein NKX3.1 |
| <u>V\$CART/CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |

| Family/matrix | Further Information |
|---------------------------------|---|
| <u>V\$NKXH/NKX25.02</u> | Homeo domain factor Nkx-2.5/Csx, tinman homolog low affinity sites |
| <u>V\$CDXF/CDX2.01</u> | Cdx-2 mammalian caudal related intestinal transcr. factor |
| <u>V\$BRNF/BRN3.01</u> | POU transcription factor Brn-3 |
| <u>V\$TBPF/TATA.02</u> | Mammalian C-type LTR TATA box |
| <u>V\$FKHD/FREAC3.01</u> | Fork head related activator-3 (FOXC1) |
| <u>V\$OCT1/OCT1.02</u> | Octamer-binding factor 1 |
| <u>V\$CART/CART1.01</u> | Cart-1 (cartilage homeoprotein 1) |
| <u>V\$PDX1/PDX1.01</u> | Pdx1 (IDX 1/IPF1) pancreatic and intestinal homeodomain TF |
| <u>V\$PARF/DBP.01</u> | Albumin D-box binding protein |
| <u>V\$GATA/GATA3.02</u> | GATA-binding factor 3 |
| <u>V\$VBPF/VBP.01</u> | PAR-type chicken vitellogenin promoter-binding protein |
| <u>V\$AP4R/TAL1ALPHA E47.01</u> | Tal-1alpha/E47 heterodimer |
| <u>V\$RP58/RP58.01</u> | Zinc finger protein RP58 (ZNF238), associated preferentially with heterochromatin |
| <u>V\$COMP/COMP1.01</u> | COMP1, cooperates with myogenic proteins in multicomponent complex |
| <u>V\$CLOX/CLOX.01</u> | Clox |
| <u>V\$TBPF/ATATA.01</u> | Avian C-type LTR TATA box |
| <u>V\$PBXC/PBX1 MEIS1.02</u> | Binding site for a Pbx1/Meis1 heterodimer |
| <u>V\$PBXF/PBX1.01</u> | Homeo domain factor Pbx-1 |

| Family/matrix | Further Information |
|------------------------|--------------------------------|
| <u>V\$IRFF/IRF1.01</u> | Interferon regulatory factor 1 |
| <u>V\$TEAF/TEF1.01</u> | TEF-1 related muscle factor |

**matches are listed in order of occurrence in the corresponding sequence

The number of consensus transcription factor binding sites present in the vector backbone (including the ampicillin resistance gene) was reduced from 224 in pGL3 to 40 in pGL4, and the number of promoter modules was reduced from 10 in pGL3 to 4 for pGL4, using databases, search programs and the like as described herein. Other modifications in pGL4 relative to pGL3 include the removal of the fl origin of replication and the redesign of the multiple cloning region.

10

MCS-1 to MCS-4 have the following sequences (SEQ ID Nos:76-79)

MCS-1

15 ACTAGTCGTCTCTCTTGAGAGACCGCGATCGCCACCATGATAAGTAA
GTAATATTAAATAAGTAAGGCCTGAGTGGCCCTCGAGCCAGCCTTGA
GTTGGTTGAGTCCAAGTCACGTCTGGAGATCTGGTACCTACGCGTGA
GCTCTACGTAGCTAGCGGCCTCGGCGGCCGAATTCTTGCGATCTAAG
TAAGCTTGGCATTCCGGTACTGTTGGTAAAGCCACCATGG

20 MCS-2

ACTAGTACGTCTCTCTTGAGAGACCGCGATCGCCACCATGATAAGTA
AGTAATATTAAATAAGTAAGGCCTGAGTGGCCCTCGAGTCAGCCTT
GAGTTGGTTGAGTCCAAGTCACGTCTGGAGATCTGGTACCTTACGCGT
AGAGCTCTACGTAGCTAGCGGCCTCGGCGGCCGAATTCTTGCGATCT
25 AAGCTTGGCAATCCGGTACTGTTGGTAAAGCCACCATGG

MCS-3

30 ACTAGTACGTCTCTCTTGAGAGACCGCGATCGCATGCCTAGGTAGGT
AGTATTAGAGCATAGGTAGAGGCCTAAGTGGCCCTCGAGTCCAGCCT
TGAGTTGGTTGAGTCCAAGTCACGTCTGGAGATCTGGTACCTTACGCG
TATGAGCTCTACGTAGCTAGCGGCCTCGGCGGCCGAATTCTTGCGAT
CTAAGCTTGGCAATCCGGTACTGTTGGTAAAGCCACCATGG

MCS-4

35 ACTAGTACGTCTCTCTTGAGAGACCGCGATCGCCACCATGTCTAGGT
AGGTAGTAAACGAAAGGGCTTAAAGGCCTAAGTGGCCCTCGAGTCCA
GCCTTGAGTTGGTTGAGTCCAAGTCACGTTTGGAGATCTGGTACCTTA

CGCGTATGAGCTCTACGTAGCTAGCGGCCTCGGCGGCCGAATTCCTTG
CGATCTAAGCTTGGCAATCCGGTACTGTTGGTAAAGCCACCATGG

bla has the following sequence:

5 ATGAGTATTCAACATTTCCGTGTCGCCCTTATTCCCTTTTTTGCGGCAT
TTTGCCTTCCTGTTTTTGTCTACCCAGAAACGCTGGTGAAAGTAAAAG
ATGCTGAAGATCAGTTGGGTGCACGAGTGGGTACATCGAACTGGAT
CTCAACAGCGGTAAGATCCTTGAGAGTTTTCGCCCCGAAGAACGTTT
TCCAATGATGAGCACTTTTAAAGTTCTGCTATGTGGCGCGGTATTATC
10 CCGTATTGACGCCGGGCAAGAGCAACTCGGTCGCCGCATACACTATT
CTCAGAATGACTTGGTTGAGTACTCACCAGTCACAGAAAAGCATCTT
ACGGATGGCATGACAGTAAGAGAATTATGCAGTGCTGCCATAACCAT
GAGTGATAAACAACGCGGCCAACTTACTTCTGACAACGATCGGAGGAC
CGAAGGAGCTAACCGCTTTTTTGCACAACATGGGGGATCATGTAAC
15 CGCCTTGATCGTTGGGAACCGGAGCTGAATGAAGCCATACCAAACGA
CGAGCGTGACACCACGATGCCTGTAGCAATGGCAACAACGTTGCGCA
AACTATTAACGCGCAACTACTTACTCTAGCTTCCCGGCAACAATTAA
TAGACTGGATGGAGGCGGATAAAGTTGCAGGACCACTTCTGCGCTCG
GCCCTTCCGGCTGGCTGGTTTTATTGCTGATAAATCTGGAGCCGGTGAG
20 CGTGGGTCTCGCGGTATCATTGCAGCACTGGGGCCAGATGGTAAGCC
CTCCCGTATCGTAGTTATCTACACGACGGGGAGTCAGGCAACTATGG
ATGAACGAAATAGACAGATCGCTGAGATAGGTGCCTCACTGATTAAG
CATTGGTAA (SEQ ID NO:41).

25 bla-1 to bla-5 have the following sequences (SEQ ID Nos:80-84):

bla-1

ACTAGTAACCCTGATAAATGCTTCAATAATATTGAAAAAGGAAGAGT
ATGAGTATTCAACATTTCCGTGTCGCCCTTATTCCCTTTTTTGCGGCAT
30 TTTGCCTTCCTGTTTTTGTCTACCCAGAAACGCTGGTGAAAGTAAAAG
ATGCTGAAGATCAGTTGGGTGCACGAGTGGGTACATCGAACTGGAT
CTCAACAGCGGTAAGATCCTTGAGAGTTTTCGCCCCGAAGAACGTTT
TCCAATGATGAGCACTTTTAAAGTTCTGCTATGTGGCGCGGTATTATC
CCGTATTGACGCCGGGCAAGAGCAACTCGGTCGCCGCATACACTATT
35 CTCAGAATGACTTGGTTGAGTACTCACCAGTCACAGAAAAGCATCTT
ACGGATGGCATGACAGTAAGAGAATTATGCAGTGCTGCCATAACCAT
GAGTGATAAACACCGCGGCCAACTTACTTCTGACAACGATCGGAGGAC
CGAAGGAGCTAACCGCTTTTTTGCACAACATGGGGGATCATGTAAC
CGCCTTGATCGTTGGGAACCGGAGCTGAATGAAGCCATACCAAACGA
40 CGAGCGTGACACCACGATGCCTGTAGCAATGGCAACAACGTTGCGCA

AACTATTAAGTGGCGAACTACTTACTCTAGCTTCCCGGCAACAATTAA
TAGACTGGATGGAGGCGGATAAAGTTGCAGGACCACTTCTGCGCTCG
GCCCTTCCGGCTGGCTGGTTTATTGCTGATAAATCTGGAGCCGGTGAG
CGTGGCTCTCGCGGTATCATTGCAGCACTGGGGCCAGATGGTAAGCC
5 CTCCCGTATCGTAGTTATCTACACGACGGGGAGTCAGGCAACTATGG
ATGAACGAAATAGACAGATCGCTGAGATAGGTGCCTCACTGATTAAG
CATTGGTAACCACTGCAGTGGTTTTCTTTTGCGGCCGC

bla-2

10 ACTAGTAACCCTGATAAATGCTGCAAACATATTGAAAAAGGAAGAGT
ATGAGTATTCAACATTTCCGTGTCGCACTCATTCCCTTCTTTGCGGCA
TTTTGCTTGCCGTGTTTTTGACACCCCCGAAACGCTGGTGAAAGTAAAA
GATGCTGAAGATCAACTGGGTGCACGAGTGGGCTATATCGAACTGGA
TCTCAATAGCGGTAAGATCCTTGAGAGTTTTTCGCCCCGAAGAACGTTT
15 TCCAATGATGAGCACTTTTAAAGTTCTGCTATGTGGCGCGGTATTATC
CCGTATTGACGCCGGGCAAGAGCAGCTCGGTCGCCGCATACACTACT
CACAGAACGACTTGGTTGAGTACTCGCCGGTCACGGAAAAGCATCTT
ACGGATGGCATGACAGTAAGAGAATTGTGTAGTGCTGCCATAACCAT
GAGTGATAAACCCGCGGCCAACTTACTTCTGACAACGATCGGAGGCC
20 CTAAGGAGCTGACCGCATTTTTTGACAACATGGGGGATCATGTAACC
CGGCTTGATCGTTGGGAACCGGAGCTGAACGAAGCCATACCGAACGA
CGAGCGTGACACCACGATGCCTGTAGCAATGGCAACAACGTTGCGCA
AACTACTCACTGGCGAACTTCTCACTCTAGCATCACGACAGCAACTC
ATAGACTGGATGGAGGCGGATAAAGTTGCAGGACCACTTCTGCGCTC
25 GGCCCTTCCGGCTGGCTGGTTTATAGCTGATAAATCCGGTGCCGGTG
AACGCGGCTCTCGCGGGATCATTGCTGCGCTGGGGCCAGATGGTAAG
CCCTCACGAATCGTAGTTATCTACACGACGGGGAGTCAGGCAACTAT
GGATGAACGAAATAGACAGATCGCTGAGATAGGTGCCTCACTGATCA
AGCACTGGTAGCCACTGCAGTGGTTTAGCTTTTGCGGCCGC

30

bla-3

ACTAGTAACCCTGACAAATGCTGCAAACATATTGAAAAAGGAAGAGT
ATGAGCATCCAACATTTTCGTGTCGCACTCATTCCCTTCTTTGCGGCA
TTTTGCTTGCCGTGTTTTTGACACCCCCGAAACGCTGGTGAAAGTAAAA
35 GATGCTGAAGATCAACTGGGTGCAAGAGTGGGCTATATCGAACTGGA
TCTCAATAGCGGCAAGATCCTTGAGTCTTTTCGCCCCGAAGAACGTTT
TCCGATGATGAGCACTTTTAAAGTTCTGCTATGTGGCGCGGTGTTGTC
CCGTATAGACGCCGGGCAAGAGCAGCTTGGTCGCCGTATACACTACT
CACAAAACGACTTGGTTGAGTACTCGCCGGTCACGGAAAAGCATCTT
40 ACGGATGGCATGACGGTAAGAGAATTGTGTAGTGCTGCCATTACCAT
GAGCGACAATAACCGCGGCCAACTTACTTCTGACAACGATCGGAGGCC
CTAAGGAGCTGACCGCATTTTTTGACAACATGGGGGATCATGTAACC
CGGCTTGACCGCTGGGAACCGGAGCTGAACGAAGCCATACCGAACG
ACGAGCGTGACACCACGATGCCTGTAGCAATGGCAACAACGTTGCGG
45 AAATACTCACTGGCGAACTTCTCACTCTAGCATCACGACAGCAGCT
CATAGACTGGATGGAGGCGGACAAAGTAGCAGGACCACTTCTTCGCT
CGGCCCTCCCTGCTGGCTGGTTTATTGCTGATAAATCCGGTGCCGGTG
AACGCGGCTCTCGCGGGATCATTGCTGCGCTGGGGCCTGATGGTAAG
CCCTCACGAATCGTAGTAATCTACACGACGGGGAGTCAGGCCACTAT

GGACGAACGAAATAGACAGATCGCTGAGATCGGTGCCTCACTGATCA
AGCACTGGTAACCACTGCAGTGGTTTAGCATTTGCGGCCGC

bla-4

5 ACTAGTAACCCTGACAAATGCTGCAAACATATTGAAAAAGGAAGAGT
ATGAGCATCCAACATTTTCGTGTCGCACTCATTCCCTTCTTTGCGGCA
TTTTGCTTGCCTGTTTTTGCACACCCCGAAACGCTGGTGAAAGTAAAA
GATGCTGAAGATCAACTGGGTGCAAGAGTGGGCTATATCGAACTGGA
10 TCTCAATAGCGGCAAGATCCTTGAGTCTTTCCGCCCCGAAGAACGTTT
TCCGATGATGAGCACTTTCAAAGTACTGCTATGTGGCGCGGTGTTGTC
CCGTATAGACGCCGGGCAAGAGCAGCTTGGTCGCCGTATACACTACT
CACAAAACGACTTGGTTGAGTACTCGCCGGTCACGGAAAAGCATCTT
ACGGATGGCATGACGGTAAGAGAATTGTGTAGTGCTGCCATTACCAT
GAGCGATAATACCGCGGCCAACTTACTTCTGACAACGATCGGAGGCC
15 CTAAGGAGCTGACCGCATTTTTTGACAACATGGGTGATCATGTGACC
CGGCTTGACCGCTGGGAACCGGAGCTGAACGAAGCCATACCGAACG
ACGAGCGTGACACCACGATGCCTGTAGCAATGGCAACAACCTCTTCGG
AAACTACTCACTGGCGAACTTCTCACTCTAGCATCACGACAGCAGCT
CATAGACTGGATGGAGGCGGACAAAGTAGCAGGACCACTTCTTCGCT
20 CGGCCCTCCCTGCTGGCTGGTTCATTGCTGATAAATCTGGAGCCGGTG
AGCGTGGCTCTCGCGGTATCATTGCTGCGCTGGGGCCTGATGGTAAG
CCCTCACGAATCGTAGTAATCTACACGACGGGGAGTCAGGCCACTAT
GGACGAACGAAATAGACAGATCGCTGAGATCGGTGCCTCACTGATCA
AGCACTGGTAACCACTGCAGTGGTTTAGCATTTGCGGCCGC

25

bla-5

ACTAGTAACCCTGACAAATGCTGCAAACATATTGAAAAAGGAAGAGT
ATGAGCATCCAACATTTTCGTGTCGCACTCATTCCCTTCTTTGCGGCA
TTTTGCTTGCCTGTTTTTGCACACCCCGAAACGCTGGTGAAAGTAAAA
30 GATGCTGAAGATCAACTGGGTGCAAGAGTGGGCTATATCGAACTGGA
TCTCAATAGCGGCAAGATCCTTGAGTCTTTCCGCCCCGAAGAACGAT
TCCCGATGATGAGCACTTTCAAAGTACTGCTATGTGGCGCGGTGTTGT
CCCGTATAGACGCCGGGCAAGAGCAGCTTGGTCGCCGTATACACTAC
TCACAAAACGACTTGGTTGAGTACTCGCCGGTCACGGAAAAGCATCT
35 TACGGATGGCATGACGGTAAGAGAATTGTGTAGTGCTGCCATTACCA
TGAGCGATAATACCGCGGCCAACTTACTTCTGACAACGATCGGAGGC
CCTAAGGAGCTGACCGCATTTTTTGACAACATGGGTGATCATGTGAC
CCGGCTTGACCGCTGGGAACCGGAGCTGAACGAAGCCATACCGAAC
GACGAGCGTGATACCACGATGCCAGTAGCAATGGCCACAACCTCTTCG
40 GAAACTACTCACTGGCGAACTTCTCACTCTAGCATCACGACAGCAGC
TCATAGACTGGATGGAGGCGGACAAAGTAGCAGGACCACTTCTTCGC
TCGGCCCTCCCTGCTGGCTGGTTCATTGCTGACAAATCCGGTGCCGGT
GAACGCGGCTCTCGCGGCATCATTGCTGCGCTGGGGCCTGATGGTAA
GCCCTCACGAATCGTAGTAATCTACACGACGGGGAGTCAGGCCACTA
45 TGGACGAACGAAATAGACAGATCGCTGAGATCGGTGCCTCACTGATC
AAGCACTGGTAACCACTGCAGTGGTTTAGCATTTGCGGCCGCNNN.

Table 30

Pairwise identity of different *bla* gene versions

| | <i>bla</i> | <i>bla-1</i> | <i>bla-2</i> | <i>bla-3</i> | <i>bla-4</i> | <i>bla-5</i> | <i>bla</i> in pGL4 (SEQ ID NO:74) |
|--------------|------------|--------------|--------------|--------------|--------------|--------------|--|
| <i>bla</i> | -- | 99 | 93 | 90 | 89 | 88 | 87 |
| <i>bla-1</i> | | -- | 94 | 90 | 90 | 89 | 88 |
| <i>bla-2</i> | | | -- | 96 | 94 | 94 | 93 |
| <i>bla-3</i> | | | | -- | 98 | 98 | 97 |
| <i>bla-4</i> | | | | | -- | 99 | 97 |
| <i>bla-5</i> | | | | | | -- | 98 |

note: sequence "bla" is bla gene from pGL3-Basic; ClustalW

(Slow/Accurate, IUB); sequence comparisons were of ORF only

5

SpeI-NcoI ver2 start has the following sequence:

ACTAGTACGTCTCTCAAGGATAAGTAAGTAATATTAAGGTACGGGAG
 GTACTTGGAGCGGCCGCAATAAAATATCTTTATTTTCATTACATCTGT
 GTGTTGGTTTTTTGTGTGAATCGATAGTACTAACATACGCTCTCCATC
 10 AAAACAAAACGAAACAAAACAACTAGCAAAATAGGCTGTCCCCAG
 TGCAAGTGCAGGTGCCAGAACATTTCTCTGGCCTAAGTGGCCGGTAC
 CGAGCTCGCTAGCCTCGAGGATATCAGATCTGGCCTCGGCGGCCAAG
 CTTGGCAATCCGGTACTGTTGGTAAAGCCACCATGG (SEQ ID NO:48);
 and

15

SpeI-NcoI-Ver2 has the following sequence:

ACTAGTACGTCTCTCAAGGATAAGTAAGTAATATTAAGGTACGGGAG
 GTATTGGACAGGCCGCAATAAAATATCTTTATTTTCATTACATCTGTG
 TGTTGGTTTTTTGTGTGAATCGATAGTACTAACATACGCTCTCCATCA
 20 AAACAAAACGAAACAAAACAACTAGCAAAATAGGCTGTCCCCAGT
 GCAAGTGCAGGTGCCAGAACATTTCTCTGGCCTAACTGGCCGGTACC
 TGAGCTCGCTAGCCTCGAGGATATCAAGATCTGGCCTCGGCGGCCAA
 GCTTGGCAATCCGGTACTGTTGGTAAAGCCACCATGG (SEQ ID NO:49)

25 pGL4 related sequences include (SEQ ID Nos.95-97):

pGL4B-4NN

GCGGCCGCAAATGCTAAACCACTGCAGTGGTTACCAGTGCTTGATCA
 30 GTGAGGCACCGATCTCAGCGATCTGTCTATTTTCGTTTCGTCCATAGTGG
 CCTGACTCCCCGTCGTGTAGATTACTACGATTCGTGAGGGCTTACCAT
 CAGGCCCCAGCGCAGCAATGATGCCGCGAGAGCCGCGTTACCGGCA

CCGGATTTGTCAGCAATGAACCAGCCAGCAGGGAGGGCCGAGCGAA
 GAAGTGGTCCTGCTACTTTGTCCGCTCCATCCAGTCTATGAGCTGCT
 GTCGTGATGCTAGAGTGAGAAGTTCGCCAGTGAGTAGTTTCCGAAGA
 GTTGTGGCCATTGCTACTGGCATCGTGGTATCACGCTCGTCGTTCCGT
 5 ATGGCTTCGTTCACTCCGTTCCAGCGGTCAAGCCGGGTCACATG
 ATCACCCATGTTGTGCAAAAATGCGGTCAGCTCCTTAGGGCCTCCGA
 TCGTTGTCAGAAGTAAGTTGGCCGCGGTATTATCGCTCATGGTAATGG
 CAGCACTACACAATTCTCTTACCGTCATGCCATCCGTAAGATGCTTTT
 CCGTGACCGGCGAGTACTCAACCAAGTCGTTTTGTGAGTAGTGTATA
 10 CGGCGACCAAGCTGCTCTTGCCCGGCGTCTATACGGGACAACACCGC
 GCCACATAGCAGTACTTTGAAAGTGCTCATCATCGGGAATCGTTCTTC
 GGGGCGGAAAGACTCAAGGATCTTGCCGCTATTGAGATCCAGTTCGA
 TATAGCCCACTCTTGACCCAGTTGATCTTCAGCATCTTTTACTTTTAC
 CAGCGTTTTCGGGGTGTGCAAAAACAGGCAAGCAAAATGCCGCAAAG
 15 AAGGGAATGAGTGCGACACGAAAATGTTGGATGCTCATACTCTTCTT
 TTTTCAATATGTTTGCAGCATTTGTCAGGGTTACTAGTACGTCTCTCTT
 GAGAGACCGCGATCGCCACCATGTCTAGGTAGGTAGTAAACGAAAG
 GGCTTAAAGGCCTAAGTGGCCCTCGAGTCCAGCCTTGAGTTGGTTGA
 GTCCAAGTCACGTTTGGAGATCTGGTACCTTACGCGTATGAGCTCTAC
 20 GTAGCTAGCGGCCTCGGCGGCCGAATTCTTGCGATCTAAGCTTGGCA
 ATCCGGTACTGTTGGTAAAGCCACCATGG

pGL4B-4NN1

gcggccgcaaatgctaaaccactgcagtggttaccagtgcctgatcagtgaggcaccgatctcagcgatctgtctatt
 25 tcgttcgtccatagtggcctgactccccgtcgttagattactacgattcgtgagggcttaccatcaggccccagcgc
 agcaatgatgccgcgagagccgcgttcaccggccccgatttgcagcaatgaaccagccagcaggaggggccg
 agcgaagaagtggctcgtactttgtccgctccatccagtctatgagctgctgctgatgtagtaagaagttc
 gccagtgagtagtttccgaagagttgtggccattgctactggcatcgtggtatcacgctcgtcgttcggtatggctcgt
 tcaactccggttccagcgggtcaagccgggtcacatgatcacccatgttgcaaaaaatcggtcagtccttaggg
 30 cctccgatcgttgcagaagtaagttggccgcggtgttcgctcatggtaatggcagcactacacaattctctaccgt
 catgccatccgtaagatgctttccgtgaccggcgagtagtcaaccaagtcgtttgtgagtagtgatacggcgacca
 agtgccttgcggcgctctatcgggacaacaccgcgcacatagcagtagtcttgaaagtgcctcatcagggaa
 tcgttctcggggcggaagactcaaggatcttgcgctattgagatccagttcgatagcccactcttcacccagt
 35 tgatcttcagcatctttactttcaccagcgttgcgggtgtgcaaaaacaggcaagcaaaatcgccgaaagaaggga
 atgagtgcgacacgaaaatgttgatgctcactcttcttttcaatatgtttgcagcatttgcagggttactagtagc
 tctctcttgagagaccgcgatcgccacatgtctaggtagtagtaaacgaaagggttaaggcctaagtggccct
 cgagtccagccttgagttggtgagtcgaagtcacgttggagatctggtaccttacgcgtatgagctctacgtagcta
 gcggcctcggcgccgaattcttgcgttcgaagcttgcaatccggtactgttgtaagccaccatgg; and

40 pGL4B-4NN2

GCGGCCGCAAATGCTAAACCACTGCAGTGGTTACCAGTGCTTGATCA
 GTGAGGCACCGATCTCAGCGATCTGCCTATTTTCGTTTCGTCCATAGTGG
 CCTGACTCCCCGTCGTGTAGATCACTACGATTCGTGAGGGCTTACCAT
 CAGGCCCCAGCGCAGCAATGATGCCGCGAGAGCCGCGTTACCCGGCC
 45 CCCGATTTGTCAGCAATGAACCAGCCAGCAGGGAGGGCCGAGCGAA
 GAAGTGGTCCTGCTACTTTGTCCGCTCCATCCAGTCTATGAGCTGCT
 GTCGTGATGCTAGAGTAAGAAGTTCGCCAGTGAGTAGTTTCCGAAGA
 GTTGTGGCCATTGCTACTGGCATCGTGGTATCACGCTCGTCGTTCCGT
 ATGGCTTCGTTCAACTCTGGTTCCAGCGGTCAAGCCGGGTCACATG

ATCACCCATGTTGTGCAAAAATGCGGTCAGCTCCTTAGGGCCTCCGA
TCGTTGTCAGAAAGTAAGTTGGCCGCGGTGTTGTCGCTCATGGTAATGG
CAGCACTACACAATTCTCTTACCGTCATGCCATCCGTAAGATGCTTTT
CCGTGACCGGCGAGTACTCAACCAAGTCGTTTTGTGAGTAGTGTATA
5 CGGCGACCAAGCTGCTCTTGCCCGGCGTCTATACGGGACAACACCGC
GCCACATAGCAGTACTTTGAAAGTGCTCATCATCGGGAATCGTTCCTC
GGGGCGGAAAGACTCAAGGATCTTGCCGCTATTGAGATCCAGTTCGA
TATAGCCCACTCTTGCACCCAGTTGATCTTCAGCATCTTTTACTTTTAC
CAGCGTTTCGGGGTGTGCAAAAACAGGCAAGCAAAATGCCGCAAAG
10 AAGGGAATGAGTGCGACACGAAAATGTTGGATGCTCATACTCTTCCT
TTTTCAATATGTTTGCAGCATTGTGTCAGGGTTACTAGTACGTCTCTCTT
GAGAGACCGCGATCGCCACCATGTCTAGGTAGGTAGTAAACGAAAG
GGCTTAAAGGCCTAAGTGGCCCTCGAGTCCAGCCTTGAGTTGGTTGA
GTCCAAGTCACGTTTGGAGATCTGGTACCTTACGCGTATGAGCTCTAC
15 GTAGCTAGCGGCCTCGGCGGCCGAATTCTTGCGTTCGAAGCTTGGCA
ATCCGGTACTGTTGGTAAAGCCACCATGG,

as well as

pGL4B-4NN3:

20 GCGGCCGCAAATGCTAAACCACTGCAGTGGTTACCAGTGCTTGATCA
GTGAGGCACCGATCTCAGCGATCTGCCTATTTTCGTTTCGTCCATAGTGG
CCTGACTCCCCGTCGTGTAGATCACTACGATTTCGTGAGGGCTTACCAT
CAGGCCCCAGCGCAGCAATGATGCCGCGAGAGCCGCGTTACCGGCC
CCCGATTTGTCAGCAATGAACCAGCCAGCAGGGAGGGCCGAGCGAA
25 GAAGTGGTCCTGCTACTTTGTCCGCCTCCATCCAGTCTATGAGCTGCT
GTCGTGATGCTAGAGTAAGAAGTTCGCCAGTGAGTAGTTTCCGAAGA
GTTGTGGCCATTGCTACTGGCATCGTGGTATCACGCTCGTCGTTCCGT
ATGGCTTCGTTCAACTCTGGTTCCAGCGGTCAAGCCGGGTCACATG
ATCACCCATATTATGAAGAAATGCAGTCAGCTCCTTAGGGCCTCCGA
30 TCGTTGTCAGAAAGTAAGTTGGCCGCGGTGTTGTCGCTCATGGTAATGG
CAGCACTACACAATTCTCTTACCGTCATGCCATCCGTAAGATGCTTTT
CCGTGACCGGCGAGTACTCAACCAAGTCGTTTTGTGAGTAGTGTATA
CGGCGACCAAGCTGCTCTTGCCCGGCGTCTATACGGGACAACACCGC
GCCACATAGCAGTACTTTGAAAGTGCTCATCATCGGGAATCGTTCCTC
35 GGGGCGGAAAGACTCAAGGATCTTGCCGCTATTGAGATCCAGTTCGA
TATAGCCCACTCTTGCACCCAGTTGATCTTCAGCATCTTTTACTTTTAC
CAGCGTTTCGGGGTGTGCAAAAACAGGCAAGCAAAATGCCGCAAAG
AAGGGAATGAGTGCGACACGAAAATGTTGGATGCTCATACTCTTCCT
TTTTCAATATGTTTGCAGCATTGTGTCAGGGTTACTAGTACGTCTCTCTT
40 GAGAGACCGCGATCGCCACCATGTCTAGGTAGGTAGTAAACGAAAG
GGCTTAAAGGCCTAAGTGGCCCTCGAGTCCAGCCTTGAGTTGGTTGA
GTCCAAGTCACGTTTGGAGATCTGGTACCTTACGCGTATGAGGGTTG
AGTCCAAGTCACGTTTGGAGATCTGGTACCTTACGCGTATGAGCTCTA
CGTAGCTAGCGGCCTCGGCGGCCGAATTCTTGCGTTCGAAGCTTGGC
45 AATCCGGTACTGTTGGTAAAGCCACCATGG (SEQ ID NO:45)

pGL4NN from Blue Heron:

GCGGCCGCAAATGCTAAACC^ACTGCAGTGGTTACCAGTGCTTGATCA

GTGAGGCACCGATCTCAGCGATCTGCCTATTTTCGTTTCGTCCATAGTGG
CCTGACTCCCGTCGTGTAGATCACTACGATTCGTGAGGGCTTACCAT
CAGGCCCCA GCGCAGCAATGATGCCGCGAGAGCCGCGTTACCGGCC
CCCGATTTGTCAGCAATGAACCAGCCAGCAGGGAGGGCCGAGCGAA
5 GAAGTGGTCCTGCTACTTTGTCCGCCTCCATCCAGTCTATGAGCTGCT
GTCGTGATGCTAGAGTAAGAAGTTCGCCAGTGAGTAGTTTCCGAAGA
GTTGTGGCCATTGCTACTGGCATCGTGGTATCACGCTCGTCGTTTCGGT
ATGGCTTCGTTCAACTCTGGTTCAGCGGTCAAGCCGGGTCACATG
ATCACCCATATTATGAAGAAATGCAGTCAGCTCCTTAGGGCCTCCGA
10 TCGTTGTCAAGTAAGTTGGCCGCGGTGTTGTCGCTCATGGTAATGG
CAGCACTACACAATTCTCTTACCGTCATGCCATCCGTAAGATGCTTTT
CCGTGACCGGCGAGTACTCAACCAAGTCGTTTTGTGAGTAGTGTATA
CGGCGACCAAGCTGCTCTTGCCCGCGTCTATACGGGACAACACCGC
GCCACATAGCAGTACTTTGAAAGTGCTCATCATCGGGAATCGTTCCTC
15 GGGGCGGAAAGACTCAAGGATCTTGCCGCTATTGAGATCCAGTTCGA
TATAGCCCACTCTTGCAACCCAGTTGATCTTCAGCATCTTTTACTTTTAC
CAGCGTTTCGGGGTGTGCAAAAACAGGCAAGCAAAATGCCGCAAAG
AAGGGAATGAGTGCGACACGAAAATGTTGGATGCTCATACTCTTCCT
TTTTCAATATGTTTGCAGCATTGTGTCAGGGTTACTAGTACGTCTCTCA
20 AGAGATTTGTGCATACACAGTGACTCATACTTTCACCAATACTTTGCA
TTTTGGATAAATACTAGACAACCTTTAGAAGTGAATTATTTATGAGGTT
GTCTTAAAAATTAAAAATTACAAAGTAATAAATCACATTGTAATGTATT
TTGTGTGATACCCAGAGGTTTAAGGCAACCTATTACTCTTATGCTCCT
GAAGTCCACAATTCACAGTCCTGAACTATAATCTTATCTTTGTGATTG
25 CTGAGCAAA TTTGCAGTATAATTTTCAGTGCTTTTAAATTTTGTCTGC
TTACTATTTTCTTTTTTATTTGGGTTTGATATGCGTGACAGAATGGG
GCTTCTATTA AAATATTCTTGAGAGACCGCGATCGCCACCATGTCTAG
GTAGGTAGTAAACGAAAGGGCTTAAAGGCCTAAGTGGCCCTCGAGTC
CAGCCTTGAGTTGGTTGAGTCCAAGTCACGTTTGGAGATCTGGTACCT
30 TACGCGTATGAGCTCTACGTAGCTAGCGGCCTCGGCGGCCGAATTCT
TGC GTTCGAAGCTTGGCAATCCGGTACTGTTGGTAAAGCCACCATGG
(SEQ ID NO:46),

pGL4 with promoter changes:

35 GCGGCCGCA AATGCTAAACCACTGCAGTGGTTACCAGTGCTTGATCA
GTGAGGCACCGATCTCAGCGATCTGCCTATTTTCGTTTCGTCCATAGTGG
CCTGACTCCCGTCGTGTAGATCACTACGATTCGTGAGGGCTTACCAT
CAGGCCCCA GCGCAGCAATGATGCCGCGAGAGCCGCGTTACCGGCC
40 CCCGATTTGTCAGCAATGAACCAGCCAGCAGGGAGGGCCGAGCGAA
GAAGTGGTCCTGCTACTTTGTCCGCCTCCATCCAGTCTATGAGCTGCT
GTCGTGATGCTAGAGTAAGAAGTTCGCCAGTGAGTAGTTTCCGAAGA
GTTGTGGCCATTGCTACTGGCATCGTGGTATCACGCTCGTCGTTTCGGT
ATGGCTTCGTTCAACTCTGGTTCAGCGGTCAAGCCGGGTCACATG
45 ATCACCCATATTATGAAGAAATGCAGTCAGCTCCTTAGGGCCTCCGA

TCGTTGTCAGAAGTAAGTTGGCCGCGGTGTTGTCGCTCATGGTAATGG
 CAGCACTACACAATTCTCTTACCGTCATGCCATCCGTAAGATGCTTTT
 CCGTGACCGGCGAGTACTCAACCAAGTCGTTTTGTGAGTAGTGTATA
 CGGCGACCAAGCTGCTCTTGCCCGGCGTCTATACGGGACAACACCGC
 5 GCCACATAGCAGTACTTTGAAAGTGCTCATCATCGGGAATCGTTCTTC
 GGGGCGGAAAGACTCAAGGATCTTGCCGCTATTGAGATCCAGTTCGA
 TATAGCCCACTCTTGCAACCAGTTGATCTTCAGCATCTTTTACTTTTAC
 CAGCGTTTCGGGGTGTGCAAAAAACAGGCAAGCAAAATGCCGCAAAG
 AAGGGAATGAGTGCGACA CGAAAATGTTGGATGCTCATACTCGTCCT
 10 TTTTCAATATTATTGAAGCATTTATCAGGGTTACTAGTACGTCTCTCA
 AGAGATTTGTGCATACACAGTGACTCATACTTTCACCAATACTTTGCA
 TTTTGGATAAATACTAGACA AACTTTAGAAGTGAATTATTTATGAGGTT
 GTCTTAAAATTAAAAATTACAAAGTAATAAATCACATTGTAATGTATT
 TTGTGTGATACCCAGAGGTTTAAGGCAACCTATTACTCTTAT (SEQ ID
 15 NO:47),

A hygromycin gene in a pGL4 vector:

Atgaagaagcccgaactcaccgctacca gcggtgaaaaatttctcatcgagaagttcgacagtgtgagcgacctgat
 20 gcagttgtcggaggcggaagagagccga gccttcagcttcgatgtcggcggacgcggctatgtactcgggtgaa
 tagctgcgctgatggcttctacaaagaccgctacgtgtaccgccacttcgccagcgctgcactacccatccccgaag
 tgttggacatcggcgagttcagcgagagccgtgacatactgcacagtagacgcgccaaggcggtactctccaaga
 cctccccgaacagagctgcctgctgtgtacagcctgtcgccgaagctatggatgctattgccgccgcgacctca
 gtcaaaccagcggttcggccattcgggc cccaaggcatcgccagtagacacaacctggcgggatttcatttgcgc
 25 cattgtgatccccatgtctaccactggcagaccgtgatggacgacaccgtgtccgccagcgtagctcaagccctgg
 acgaactgatgctgtgggccaagactgtccgaggtgcgccacctgtccatgccacttcggcagcaacaacgt
 cctgaccgacaacggccgcatcaccgcc gtaatcgactgttccgaagctatgttcggggacagtcagtagaggtg
 gccaacatcttcttggcgccctggctg gcttgcattggagcagcagactcgtacttcgagcgccggcatccga
 gctggccggcagccctcgtctgcgagcctacatgctgcgcatcggcctggatcagctctaccagagcctcgtggac
 30 ggcaacttcgacgatgctgcctgggtcaaaggccgtgcgatgccatcgtccgcagcggggccggcaccgtcggt
 cgcacacaaatcgctcgggagcgcagccgtatggaccgacggctgcgtcgaggtgctggccgacagcggca
 accgcccggccagtagacgaccgcgcgctaaggaggtaggtcgagtttaa (SEQ ID NO:88),

35 pGL4.10

ggcctaactggccggtacctgagctcgcta gcctcgaggatatcaagatctggcctcggcggccaagcttggcaat
 ccggtactgttgtaaagccaccatggaagatgccaaaaacattaagaaggggccagcgccattctaccactcga
 agacgggaccgcccggcgagcagctgcacaaagccatgaagcgctacgccctggtgcccgccaccatcgcttta
 40 ccgacgcacatacaggtggacattacctacgccgagtacttcgagatgagcgttcggctggcagaagctatgaa
 gcgctatgggctgaatacaaacatcggtcgtggtgtgcagcgagaatagcttcagttctcatgccgtgttggg
 tggcctgttcatcggtgtggctgtggccca gctaacgacatctacaacgagcgcgagctgctgaacagcatgggc

atcagccagcccaccgtcgtattcgtgagcaagaaagggtgcaaaagatcctcaacgtgcaaaagaagctaccg
 atcataaaaagatcatcatcatgtagcaagaccgactaccagggttc.aaaagcatgtacaccttcgtgacttcc
 catttgcacccgggttaacaggtacgacttcgtgccgagagcttcgac.aggacaaaaccatcgccctgatcat
 gaacagtagtggcagtagccgattgccaagggtgtagccctaccgcac.aggaccgttgtgtccgattcagtcac
 5 gcccgcgaccccatcttcggcaaccagatcatccccgacaccgctatcc.aggcgtgggtccatttcaccacggctt
 cggcatgttcaccacgtgggtacttgatctgcggcttccgggtcgtgctc.atgtaccgttcgaggaggagctattc
 ttgcgcagcttgcaagactataagattcaatctgccctgctgggtgccacac.tatttagcttcttcgtaagagcactct
 catcgacaagtacgacctaagcaacttcacgagatcgccagcggcgggg.gcgccgctcagcaaggaggtaggtg
 aggccgtggccaaacgcttcacctaccaggcatccgacgggtacgg.cctgacagaacaaccagcgccattc
 10 tcatcccccgaaagggtgacgacaagcctggcgagtaggcaagggtgtgccccttctc.gaggctaagggtgtgtg
 acttggacaccggtaagacactgggtgtgaaccagcggcgagctgtg.cgtccgtggcccatgatcatgagcg
 gtacgttaacaacccgaggtacaaacgctctcatcgacaaggacggc.tggctgcacagcggcgacatcgcc
 actgggacgaggacgagcacttctcatctggaccgggtgaagagcctg.atcaatacaagggtaccaggtagc
 cccagccgaactggagagcatcctgctgaacaccccaacatcttcgacg.cgggggtcgccggcctgcccgcag
 15 acgatccggcgagctgcccggcgagctgctgctggaacacggtaaaaccatgaccgagaaggagatcgtg
 gactatgtggccagccaggttacaaccgccaagaagctgcgcgggtgtgt.tgtgttcgtggacgaggtgcctaaag
 gactgaccggcaagttggacgcccgaagatccgcgagattctattaag.gccaagaagggtggcaagatcgcc
 gtgtaataattctagatcgggggcgccggcgcttcgagcagacatgata.agatacattgatgagttggacaaac
 cacaactagaatgcagtgaaaaaatgcttatttgtgaaatttgtatgctat.tgcttatttgaaccattataagctgca
 20 ataaacaagttaacaacaacaattgcattcattttatgtttcagggttcaggggg.agggtgtgggaggtttttaaagcaagt
 aaaacctctacaatatgtgtaaaatcgataaggatccgtcgaccgatgccc.ttgagagcctcaaccagtcagctcc
 ttccgggtgggcgcggggcatgactatcgctcgccgacttatgactgtctct.ttatcatgcaactcgtaggacaggtgc
 cggcagcgctcttccgcttctcgtcactgactcgtcgctcggtcgttc.ggctgcggcgagcggtatcagctca
 ctcaaaggcggtataacggttatccacagaatcaggggataacgcaggaa.agaacatgtgagcaaaaggccagca
 25 aaaggccaggaaccgtaaaaaggccgctgtgctggcgttttccataggct.cggccccctgacgagcatcacaaa
 aatcgacgctcaagttagaggtggcgaaacccgacaggactataagata.ccaggcggttccccctggaagctccc
 tctgctcgtctctgttccgacctgcccgttaccggatacctgtccgccttt.ctcccttcgggaagcggtggcgttct
 catagctcacgctgtaggtatctcagttcggtgtaggtcgctcgctcaagct.gggctgtgtgcacgaacccccgttc
 agcccgaccgctcgccctatccgtaactatcgcttgagttcaacccggt.aagacacgacttatcgccactggca
 30 gcagccactggttaacaggattagcagagcgaggtatgtaggcgggtgtac.agagttctgaagtgtgacctaaacta
 cggctacactagaagaacagatttggatctgcgctctgctgaagccagttacc.ttcggaaaaagagttgtagctct
 tgatccggcaaaacaaaccaccgctgtagcggtgtgtttttgttgaagca.gcagattacgcgcagaaaaaaagg
 atctcaagaagatccttgatctttctacgggtctgacgctcagtggaacg.aaactcacgtaagggtttgtgca
 tgagattatcaaaaaggatctcacctagatccttttaataaaaaatgaagttt.taaatcaatcaaatatataatgagta
 35 aacttggctgtacagcggccgcaaatgctaaaccactcgagtggttaccagt.gcttgatcagtgaggcaccgatctc
 agcgatctgcctatttctgttcgcatagtgccctgactccccgctgtagat.cactacgattcgtgagggttaccat
 caggccccagcgcaaatgatccgcgagagccggttaccggcccc.cgatttgcagcaaatgaaccagcca
 gcaggggagggccgagcgaagaagtggctcctgctactttgtccgcctccatc.cagtctatgagctgctgctgatgc
 tagagtaagaagttcgccagttagtagttccgaagagttgtggccattgct.actggcatcgtggtatcacgctcgtc
 40 ttgggtatggcttctgtaactctggttccagcggtcaaggcggtcacatg.atcaccataattatgaagaatgcag
 tcagctccttagggcctccgctgttcagaagtaagttggccggtgtgtgctcgtcatggtaatggcagcactac
 acaattctcttaccgtcatgccatccgtaagatgcttttccgtgaccggcgag.tactcaaccaagtcgttttgtgagtagt
 gtatacggcgaccaagctgcttgcggcgctctatcgggacaacaccg.cgccacatagcagtagcttgaagtgtg
 ctcatcatcggaatcgttcttcggggcggaagactcaaggatcttgcgc.tattgagatccagttcgatatagccc
 45 actcttcacccagttgatcttcagcatcttttacttccaccagcggttcgggtgtgcaaaaacaggcaagcaaaatgc
 cgcaagaagggaatgagtgcgacacgaaaatgttgatgctcatactgt.ccttttcaatattatgaagcatttate
 aggggtactagtagctctcaaggataagtaaatattaaggtagggga.ggtattggacagggccgcaataaaata
 tctttatttcttaccatctgtgtgtgtgttttgtgtgaatcgatagtagtaacat.cgcctctccatcaaaaacaaacgaaa
 caaaacaaactagcaaaatagggtgtccccagtgcaagtgaggtgccagaacatttctctaagtaataatgaagtagt

catacgctctccatcaaaacaaaacgaacaaaacaaactagcaaaaataggctgtccccagtgcaagtcaggtgc
cagaacatttctct (SEQ ID NO:90).

The pGL4 backbone (*NotI*-*NcoI*) has the following sequence:

5 gcggccgcaaatgctaaaccactgcagtggttaccagtgcctgatcagtgaggcaccgatctcagcgatctgcctatt
tcgttcgtccatagtgccctgactccccgtcgtgtagatcactacgattcgtgagggcttaccatcagggcccgccg
agcaatgatgccgcgagagccgcgttcaccggcccccgattgtcagcaatgaaccagccagcagggaggccg
agcgaagaagtggctcctgtactttgtccgcctccatccagtctatgagctgctgctgtagctagtaagaagttc
gccagtgagtagttccgaagagttgtggccattgctactggcatcgtggtatcacgctcgtcgttcggtatggcttgc
10 tcaactctggttcccagcggtcaagccgggtcacatgatcacccatattatgaagaaatgcagtcagctccttagggc
ctccgatcgttgcagaagtaagttggccgggtgtgtcgtcatggtaatggcagcactacacaattctctaccgtc
atgcatccgtaagatgctttccgtgaccggcgagctactcaaccaagtcgtttgtgagtagtgatatacggcgaccaa
gctgctcttgcggcgtctatacgggacaacaccgcgccacatagcagctactttgaaagtgtcatcatcggggaat
cgttcttcggggcgaaagactcaagatcttgcgctattgagatccagttcgatagcccactcttgacccagtt
15 gatcttcagcatcttttactttaccagcggttcggggtgtgcaaaaacaggcaagcaaaatccgcgaagaaggga
atgagtgcgacacgaaaatgttgatgctcactcgtccttttcaatattattgaagcatttatcagggttactagtacg
tcttcaaggataagtaagtaataattaaggtagggaggtattggacaggccgcaataaaatactttattttcattacat
ctgtgtgttggtttttgtgtaatcgatagtaactaacatagctctccatcaaaacaaaacgaaacaaaacaaactagc
aaaataggctgtccccagtgcaagtcaggtgccagaacatttcttgccctaactggccggtacctgagctcgcta
20 gcctcgaggatatcaagatctggcctcggcgccaagcttggaatccggtactgttgtaagccaccatgg
(SEQ ID NO:74).

Example 10

Summary of Sequences Removed in Synthetic Genes

25 Search parameters:

TFBS searches were limited to vertebrate TF binding sites. Searches were performed by matrix family, *i.e.*, the results show only the best match from a family for each site. MatInspector default parameters were used for the core and matrix similarity values (core similarity = 0.75, matrix similarity = optimized), except for sequence MCS-1 (core similarity = 1.00, matrix similarity = optimized).

Promoter module searches included all available promoter modules (vertebrate and others) and were performed using default parameters (optimized threshold or 80% of maximum score).

35 Splice site searches were performed for splice acceptor or donor consensus sequences.

Table 31

| Sequence | Matrix Library | TFBS (family matches) | Promoter modules | Splice sites (+ strand) |
|-------------|--------------------|-----------------------|------------------|-------------------------|
| puro | (not applicable) | 62 | 5 | 0 |
| hpuro | (not applicable) | 68 | 4 | 1 |
| hpuro1 | Ver 4.1 Feb 2004 | 4 | 2 | 1 |
| hpuro2 | Ver 4.1 Feb 2004 | 2 | 0 | 1 |
| ----- | ----- | ----- | ----- | ----- |
| Neo | (not applicable) | 53 | 0 | No data |
| hneo | (not applicable) | 61 | 2 | 3 |
| hneo-1 | Ver 3.1.2 Jun 2003 | No data | No data | No data |
| hneo-2 | Ver 3.1.2 Jun 2003 | No data | No data | No data |
| hneo-3 | Ver 3.1.2 Jun 2003 | 0 | 0 | 0 |
| hneo-4 | Ver 4.1 Feb 2004 | 7 | 1 | 0 |
| hneo-5 | Ver 4.1 Feb 2004 | 0 | 0 | 0 |
| ----- | ----- | ----- | ----- | ----- |
| Hyg | (not applicable) | 74 | 3 | No data |
| hhyg | (not applicable) | 94 | 4 | 6 |
| hhyg-1 | Ver 3.1.2 Jun 2003 | No data | No data | No data |
| hhyg-2 | Ver 3.1.2 Jun 2003 | No data | No data | No data |
| hhyg-3 | Ver 3.1.2 Jun 2003 | 3 | 0 | 0 |
| hHygro | Ver 3.3 Aug 2003 | 5 | 0 | 0 |
| hhyg-4 | Ver 3.3 Aug 2003 | 4 | 0 | 0 |
| ----- | ----- | ----- | ----- | ----- |
| Luc | (not applicable) | 213 | 11 | No data |
| Luc+ | (not applicable) | 189 | 7 | No data |

| Sequence | Matrix Library | TFBS (family matches) | Promoter modules | Splice sites (+ strand) |
|--------------|--------------------|-----------------------|------------------|-------------------------|
| | applicable) | | | |
| hluc+ver2A1 | Ver 3.0 Nov 2002 | 110 | 7 | 6 |
| hluc+ver2A2 | Ver 3.0 Nov 2002 | No data | No data | No data |
| hluc+ver2A3 | Ver 3.0 Nov 2002 | 8 | No data | 0 |
| hluc+ver2A4 | Ver 3.0 Nov 2002 | No data | No data | No data |
| hluc+ver2A5 | Ver 3.0 Nov 2002 | No data | No data | No data |
| hluc+ver2A6 | Ver 3.0 Nov 2002 | 2 | 0 | 0 |
| hluc+ver2A6 | Ver 3.1.1 Apr 2003 | 4 | 0 | 0 |
| hluc+ver2A7 | Ver 3.1.1 Apr 2003 | 1 | 0 | 0 |
| hluc+ver2A8 | Ver 3.1.1 Apr 2003 | 1 | 0 | 0 |
| hluc+ver2B1 | Ver 3.0 Nov 2002 | 187 | 2 | 8 |
| hluc+ver2B2 | Ver 3.0 Nov 2002 | No data | No data | No data |
| hluc+ver2B3 | Ver 3.0 Nov 2002 | 35 | No data | 0 |
| hluc+ver2B4 | Ver 3.0 Nov 2002 | No data | No data | No data |
| hluc+ver2B5 | Ver 3.0 Nov 2002 | No data | No data | No data |
| hluc+ver2B6 | Ver 3.0 Nov 2002 | 2 | 0 | 0 |
| hluc+ver2B6 | Ver 3.1.1 Apr 2003 | 6 | 0 | 0 |
| hluc+ver2B7 | Ver 3.1.1 Apr 2003 | 2 | 0 | 0 |
| hluc+ver2B8 | Ver 3.1.1 Apr 2003 | 1 | 0 | 0 |
| hluc+ver2B9 | Ver 3.1.1 Apr 2003 | 1 | 0 | 0 |
| hluc+ver2B10 | Ver 3.1.1 Apr 2003 | 1 | 0 | 0 |
| ----- | ----- | ----- | ----- | ----- |
| MCS-1 | Ver 2.2 Sep 2001 | 14 | No data | (not applicable) |

| Sequence | Matrix Library | TFBS (family matches) | Promoter modules | Splice sites (+ strand) |
|-----------------------------|------------------|-----------------------|------------------|-------------------------|
| MCS-2 | Ver 2.2 Sep 2001 | 12 | No data | (not applicable) |
| MCS-3 | Ver 2.2 Sep 2001 | 0 | No data | (not applicable) |
| MCS-4 | Ver 2.3 Feb 2001 | 0 | 0 | (not applicable) |
| ----- | ----- | ----- | ----- | ----- |
| Bla | (not applicable) | No data | No data | (not applicable) |
| bla-1 | Ver 2.2 Sep 2001 | 94 | 1 | (not applicable) |
| bla-2 | Ver 2.3 Feb 2001 | 51 | No data | (not applicable) |
| bla-3 | Ver 2.3 Feb 2001 | 16 | No data | (not applicable) |
| bla-4 | Ver 2.3 Feb 2001 | 14 | No data | (not applicable) |
| bla-5 | Ver 2.3 Feb 2001 | 5 | 0 | (not applicable) |
| ----- | ----- | ----- | ----- | ----- |
| pGL4B-4NN | Ver 2.4 May 2002 | 11 | 0 | (not applicable) |
| pGL4B-4NN1 | Ver 2.4 May 2002 | 7 | No data | (not applicable) |
| pGL4B-4NN2 | Ver 2.4 May 2002 | 4 | 0 | (not applicable) |
| pGL4B-4NN3 | Ver 2.4 May 2002 | 3 | 0 | (not applicable) |
| ----- | ----- | ----- | ----- | ----- |
| <i>SpeI-NcoI-Ver2-Start</i> | Ver 4.0 Nov 2003 | 34 | 1 | (not applicable) |
| <i>SpeI-NcoI-Ver2</i> | Ver 4.0 Nov 2003 | 28 | 1 | (not applicable) |

Using the 5 sequences, i.e., hluc+ver2A1, bla-1, hneo-1, hpuro-1, hhyg-1 (humanized codon usage) for analysis, TFBS from the following families were found in 3 out 5 sequences:

- 5 V\$AHRR (AHR-arnt heterodimers and AHR-related factors)
 V\$ETSF (Human and murine ETS1 factors)
 V&NFKB (Nuclear Factor Kappa B/c-rel)

- V\$VMYB (AMV-viral myb oncogene)
- V\$CDEF (Cell cycle regulators: Cell cycle dependent element)
- V\$HAND (bHLH transcription factor dimer of HAND2 and E12)
- V\$NRSF (Neuron-Restrictive Silencer Factor)
- 5 V\$WHZF (Winged Helix and ZF5 binding sites)
- V\$CMYB (C-myb, cellular transcriptional activator)
- V\$MINI (Muscle INitiator)
- V\$P53F (p53 tumor suppr.-neg. regulat. of the tumor suppr. Rb)
- V\$ZF5F (ZF5 POZ domain zinc finger)
- 10 V\$DEAF (Homolog to deformed epidermal autoregulatory factor-1
from *D. melanogaster*)
- V\$MYOD (MYOblast Determining factor)
- V\$PAX5 (PAX-5/PAX-9 B-cell-specific activating protein)
- V\$EGRF (EGR/nerve growth Factor Induced protein C & rel. fact.)
- 15 V\$NEUR (NeuroD, Beta2, HLH domain)
- V\$REBV (Epstein-Barr virus transcription factor R);
- TFBS from the following families were found in 4 out of 5 sequences:
- V\$ETSF (Human and murine ETS1 factors)
- 20 V\$CDEF (Cell cycle regulators: Cell cycle dependent element)
- V\$HAND (bHLH transcription factor dimer of HAND2 and E12)
- V\$NRSF (Neuron-Restrictive Silencer Factor)
- V\$PAX5 (PAX-5/PAX-9 B-cell-specific activating protein)
- V\$NEUR (NeuroD, Beta2, HLH domain); and
- 25 TFBS from the following families were found in 5 out of 5 sequences:
- V\$PAX5 (PAX-5/PAX-9 B-cell-specific activating protein).

30 References

- Altschul et al., Nucl. Acids Res., 25, 3389 (1997).
- Aota et al., Nucl. Acids Res., 16, 315 (1988).
- Boshart et al., Cell, 41, 521 (1985).

- Bronstein et al., Cal. Biochem., 219, 169 (1994).
- Corpet et al., Nucl. Acids Res., 16, 881 (1988).
- deWet et al., Mol. Cell. Biol., 7, 725 (1987).
- Dijkema et al., EMBO J., 4, 761 (1985).
- 5 Faist and Meyer, Nucl. Acids Res., 20, 26 (1992).
- Gorman et al., Proc. Natl. Acad. Sci. USA, 79, 6777 (1982).
- Higgins et al., Gene, 73, 237 (1985).
- Higgins et al., CABIOS, 5, 151 (1989).
- Huang et al., CABIOS, 8, 155 (1992).
- 10 Itolcik et al., PNAS, 94, 12410 (1997).
- Johnson et al., Mol. Reprod. Devel., 50, 377 (1998).
- Jones et al., Mol. Cell. Biol., 17, 6970 (1997).
- Karlin and Altschul, Proc. Natl. Acad. Sci. USA, 87, 2264 (1990).
- Karlin and Altschul, Proc. Natl. Acad. Sci. USA, 90, 5873 (1993).
- 15 Keller et al., J. Cell Biol., 84, 3264 (1987).
- Kim et al., Gene, 91, 217 (1990).
- Lamb et al., Mol. Reprod. Devel., 51, 218 (1998).
- Mariatis et al., Science, 236, 1237 (1987).
- Michael et al., EMBO. J., 9, 481 (1990).
- 20 Mizushima and Nagata, Nucl. Acids Res., 18, 5322 (1990).
- Murray et al., Nucl. Acids Res., 17, 477 (1989).
- Myers and Miller, CABIOS, 4, 11 (1988).
- Nakamura et al., NAR, 28:292 (2000).
- Needleman and Wunsen, J. Mol. Biol., 48, 443 (1970).
- 25 Pearson and Lipman, Proc. Natl. Acad. Sci. USA, 85, 2444 (1988).
- Pearson et al., Meth. Mol. Biol., 24, 307 (1994).
- Sharp et al., Nucl. Acids Res., 16, 8207 (1988).
- Sharp et al., Nucl. Acids Res., 15, 1281 (1987).
- Smith and Waterman, Adv. Appl. Math., 2, 482 (1981).
- 30 Stemmer et al., Gene, 164, 49 (1995).
- Uetsuki et al., J. Biol. Chem., 264, 5791 (1989).
- Voss et al., Trends Biochem. Sci., 11, 287 (1986).
- Wada et al., Nucl. Acids Res., 18, 2367 (1990).

Watson et al, eds. Recombinant DNA: A Short Course, Scientific American Books, W. H. Freeman and Company, New York (1983).

Wood, K. Photochemistry and Photobiology, 62, 662 (1995).

Wood, K. Science 244, 700 (1989)

5

All publications, patents and patent applications are incorporated herein by reference. While in the foregoing specification, this invention has been described in relation to certain preferred embodiments thereof, and many details have been set forth for purposes of illustration, it will be apparent to those skilled in the art that the invention is susceptible to additional embodiments and that certain of the details herein may be varied considerably without departing from the basic principles of the invention.

10

WHAT IS CLAIMED IS:

1. An isolated nucleic acid molecule comprising a synthetic nucleotide sequence having a coding region for a selectable polypeptide, wherein the synthetic nucleotide sequence has 90% or less nucleic acid sequence identity to a parent nucleic acid sequence encoding a corresponding selectable polypeptide, wherein the decreased sequence identity is a result of different codons in the synthetic nucleotide sequence relative to the codons in the parent nucleic acid sequence, wherein the nucleotide sequence encodes a selectable polypeptide with at least 85% amino acid sequence identity to the corresponding selectable polypeptide encoded by the parent nucleic acid sequence, wherein the synthetic nucleotide sequence has a reduced number of regulatory sequences relative to the average number of regulatory sequences resulting from random selections of codons at the sequences which differ between the synthetic nucleotide sequence and the parent nucleic acid sequence, and wherein the synthetic nucleotide sequence, when expressed in a cell, confers resistance to ampicillin, puromycin, hygromycin or neomycin.
2. The isolated nucleic acid molecule of claim 1 wherein the regulatory sequences include transcription factor binding sequences, intron splice sites, poly(A) sites, promoter modules, and/or promoter sequences.
3. The isolated nucleic acid molecule of claim 1 wherein a majority of the codons which differ are ones that are preferred codons of a desired host cell and/or are not low-usage codons in that host cell.
4. The isolated nucleic acid molecule of claim 3 wherein the majority of the codons which differ in the synthetic nucleic acid sequence are those which are employed more frequently in mammals.

5. The isolated nucleic acid molecule of claim 3 wherein the majority of the codons which differ in the synthetic nucleic acid sequence are those which are preferred codons in humans.
- 5 6. The isolated nucleic acid molecule of claim 3 wherein the majority of codons which differ are the codons CGC, CTG, AGC, ACC, CCC, GCC, GGC, GTG, ATC, AAG, AAC, CAG, CAC, GAG, GAC, TAC, TGC and TTC.
- 10 7. The isolated nucleic acid molecule of claim 1 wherein the nucleic acid molecule encodes a fusion of the selectable polypeptide with a luciferase.
8. The isolated nucleic acid molecule of claim 7 wherein the luciferase is a *Renilla* luciferase, a firefly luciferase or a click beetle luciferase.
- 15 9. The isolated nucleic acid molecule of claim 1 wherein the parent nucleic acid sequence is a wild-type *neo*, *hyg*, *bLa* or *puro* sequence.
- 10 20 10. The isolated nucleic acid molecule of claim 1 wherein the parent nucleic acid sequence is SEQ ID NO:1, SEQ ID NO:6, SEQ ID NO:15 or SEQ ID NO:41.
- 25 11. The isolated nucleic acid molecule of claim 1 wherein the synthetic nucleotide sequence comprises an open reading frame in SEQ ID NO:4, SEQ ID NO:5, SEQ ID NO:9, SEQ ID NO:10, SEQ ID NO:11, SEQ ID NO:30, SEQ ID NO:38, SEQ ID NO:39, SEQ ID NO:42, SEQ ID NO:44; SEQ ID NO:70, SEQ ID NO:71, SEQ ID NO:72, SEQ ID NO:73, SEQ ID NO:74, SEQ ID NO:80, SEQ ID NO:81, SEQ ID NO:82, SEQ ID NO:83, or SEQ ID NO:84.
- 30 12. The isolated nucleic acid molecule of claim 1 wherein the synthetic nucleotide sequence has at least 10% fewer regulatory sequences.

13. The isolated nucleic acid molecule of claim 1 wherein the synthetic nucleotide sequence has an increased number of AGC serine-encoding codons, an increased number of ATC isoleucine-encoding codons, an increased number of CCC proline-encoding codons, and/or an increased number of ACC threonine-encoding codons.
14. The isolated nucleic acid molecule of claim 1 wherein the codons in the synthetic nucleotide sequence which differ encode the same amino acids as the corresponding codons in the parent nucleic acid sequence.
15. The isolated nucleic acid molecule of claim 1 which has at least 90% nucleotide sequence identity to an open reading frame in any one of SEQ ID NO:4, SEQ ID NO:5, SEQ ID NO:9, SEQ ID NO:10, SEQ ID NO:11, SEQ ID NO:30, SEQ ID NO:38, SEQ ID NO:39, SEQ ID NO:42, SEQ ID NO:44, SEQ ID NO:70, SEQ ID NO:71, SEQ ID NO:72, SEQ ID NO:73, SEQ ID NO:74, SEQ ID NO:80, SEQ ID NO:81, SEQ ID NO:82, SEQ ID NO:83, or SEQ ID NO:84, or the complement thereof.
16. The isolated nucleic acid molecule of claim 1 wherein the nucleic acid molecule encodes a fusion of the selectable polypeptide with one or more other peptides or polypeptides, wherein at least the selectable polypeptide is encoded by the synthetic nucleic acid sequence.
17. The isolated nucleic acid molecule of claim 16 wherein one or more other peptides are peptides having protein destabilization sequences.
18. A plasmid comprising the nucleic acid molecule of claim 1.
19. The plasmid of claim 18 which further comprises a multiple cloning region.
20. The plasmid of claim 18 which further comprises an open reading frame of interest.

21. The plasmid of claim 18 which further comprises a promoter functional in a particular host cell operably linked to the synthetic nucleotide sequence.
- 5 22. The plasmid of claim 21 wherein the promoter is functional in a prokaryotic cell.
23. The plasmid of claim 21 wherein the promoter is functional in a eukaryotic cell.
- 10 24. The plasmid of claim 20 further comprising a promoter operably linked to the open reading frame of interest.
- 15 25. An isolated nucleic acid molecule comprising a synthetic nucleotide sequence encoding a firefly luciferase, wherein the synthetic nucleotide sequence has 80% or less nucleic acid sequence identity to a parent nucleic acid sequence having SEQ ID NO:43 or 85% or less nucleic acid sequence identity to a parent nucleic acid sequence having SEQ ID
- 20 NO:14 which encodes a firefly luciferase, wherein the decreased sequence identity is a result of different codons in the synthetic nucleotide sequence relative to the codons in the parent nucleic acid sequence, wherein the synthetic nucleotide sequence encodes a firefly luciferase which has at least 85% amino acid sequence identity to the
- 25 corresponding luciferase encoded by the parent nucleic acid sequence, and wherein the synthetic nucleotide sequence has a reduced number of regulatory sequences relative to the average number of regulatory sequences resulting from random selections of codons at the sequences which differ between the synthetic nucleotide sequence and the parent
- 30 nucleic acid sequence.
26. The isolated nucleic acid molecule of claim 25 wherein the regulatory sequences include transcription factor binding sequences, intron splice

sites, poly(A) sites, promoter modules, and/or promoter sequences.

27. The isolated nucleic acid molecule of claim 25 wherein a majority of the codons which differ are ones that are preferred codons of a desired host cell and/or are not low-usage codons in that host cell.
28. The isolated nucleic acid molecule of claim 27 wherein the majority of the codons which differ in the synthetic nucleic acid molecule are those which are employed more frequently in mammals.
29. The isolated nucleic acid molecule of claim 27 wherein the majority of the codons which differ in the synthetic nucleic acid molecule are those which are preferred codons in humans.
30. The isolated nucleic acid molecule of claim 27 wherein the majority of codons which differ are the codons CGC, CTG, AGC, ACC, CCC, GCC, GGC, GTG, ATC, AAG, AAC, CAG, CAC, GAG, GAC, TAC, TGC and TTC.
31. The isolated nucleic acid molecule of claim 25 wherein the synthetic nucleotide sequence comprises a sequence in an open reading frame in SEQ ID NO:21, SEQ ID NO:22, or SEQ ID NO:23 or has at least 90% nucleotide sequence identity thereto.
32. The isolated nucleic acid molecule of claim 25 wherein the synthetic nucleic acid molecule is expressed in a mammalian host cell at a level which is greater than that of the parent nucleic acid sequence.
33. The isolated nucleic acid molecule of claim 25 wherein the synthetic nucleic acid molecule has an increased number of AGC serine-encoding codons, an increased number of CCC proline-encoding codons, an increased number of ATC isoleucine-encoding codons and/or an increased number of ACC threonine-encoding codons.

34. The isolated acid molecule of claim 25 wherein the synthetic nucleotide sequence has at least 10% fewer transcription regulatory sequences .
- 5 35. The isolated nucleic acid molecule of claim 25 wherein the codons in the synthetic nucleotide sequence which differ encode the same amino acids as the corresponding codons in the parent nucleic acid sequence.
- 10 36. The isolated nucleic acid molecule of claim 25 wherein the nucleic acid molecule encodes a fusion of the luciferase with one or more other peptides or polypeptides, wherein at least the luciferase is encoded by the synthetic nucleic acid sequence.
- 15 37. The isolated nucleic acid molecule of claim 36 wherein one or more other peptides are peptides having protein destabilization sequences.
38. A plasmid comprising the nucleic acid molecule of claim 25.
- 20 39. The plasmid of claim 38 which further comprises a multiple cloning region.
40. The plasmid of claim 38 which further comprises a promoter operatively linked to the synthetic nucleotide sequence.
- 25 41. The plasmid of claim 38 which further comprises the synthetic nucleotide sequence of the nucleic acid molecule of claim 1.
42. An expression vector comprising the nucleic acid molecule of claim 25 linked to a promoter functional in a cell.
- 30 43. The expression vector of claim 42 wherein the promoter is functional in a eukaryotic cell.

44. The expression vector of claim 42 wherein the expression vector further comprises a multiple cloning site.
45. The expression vector of claim 42 wherein the promoter is functional in a mammalian cell.
46. The expression vector of claim 42 wherein the synthetic nucleotide sequence is operatively linked to a Kozak consensus sequence.
47. A plasmid comprising a nucleotide sequence comprising SEQ ID NO:74 or a nucleotide sequence comprising at least 80% nucleic acid sequence identity to SEQ ID NO:74, which nucleotide sequence comprises an open reading frame with less than 90% nucleic acid sequence identity to SEQ ID NO:41, and the expression of which open reading frame in a host cell confers resistance to ampicillin.
48. A host cell comprising the expression cassette of claim 42.
49. A host cell comprising the plasmid of claim 17, 38 or 47.
50. A kit comprising, in suitable container means, the plasmid of claim 17, 38 or 47.
51. A polynucleotide which hybridizes under stringent hybridization conditions to SEQ ID NO:4, SEQ ID NO:5, SEQ ID NO:9, SEQ ID NO:10, SEQ ID NO:11, SEQ ID NO:30, SEQ ID NO:38, SEQ ID NO:39, SEQ ID NO:42, SEQ ID NO:44, SEQ ID NO:70, SEQ ID NO:71, SEQ ID NO:72, SEQ ID NO:73, SEQ ID NO:74, SEQ ID NO:80, SEQ ID NO:81, SEQ ID NO:82, SEQ ID NO:83, SEQ ID NO:84, SEQ ID NO:21, SEQ ID NO:22, SEQ ID NO:23, or the complement of the polynucleotide, wherein the polynucleotide or the complement thereof encodes a selectable polypeptide or a firefly luciferase.

52. The polynucleotide of claim 51 which does not have SEQ ID NO:1, SEQ ID NO:6, SEQ ID NO:15, SEQ ID NO:41, SEQ ID NO:14, or SEQ ID NO:43.
- 5
53. An isolated nucleic acid molecule comprising a synthetic nucleotide sequence which does not code for a desirable peptide or polypeptide but includes sequences which inhibit transcription and/or translation, wherein the synthetic nucleotide sequence has at least 20 nucleotides which have
- 10 a different sequence relative to a corresponding parent nucleic acid sequence which does not code for the desirable peptide or polypeptide, wherein the synthetic nucleotide sequences has 90% or less nucleic acid sequence identity to the parent nucleic acid sequences, and wherein the sequence difference is a result of a reduced number of one or more
- 15 regulatory sequences in the synthetic nucleotide sequence relative to the parent nucleic acid sequence.
54. The isolated nucleic acid molecule of claim 53 wherein the synthetic nucleotide sequence has SEQ ID NO:49.
- 20
55. The isolated nucleic acid molecule of claim 53 further comprising a multiple cloning region and/or a poly(A) site.
56. The isolated nucleic acid molecule of claim 53 wherein the sequences
- 25 which inhibit transcription include one or more poly(A) sites.
57. The isolated nucleic acid molecule of claim 53 wherein the sequences which inhibit translation include one or more stop codons in one or more reading frames.
- 30
58. The isolated nucleic acid molecule of claim 53 wherein the parent nucleic acid sequence includes a multiple cloning region.

59. The isolated nucleic acid molecule of claim 53 wherein the parent nucleic acid sequence includes sequences which inhibit transcription and/or translation.
- 5 60. The isolated nucleic acid molecule of claim 53 wherein the parent nucleic acid sequence has SEQ ID NO:76.
61. The isolated nucleic acid molecule of claim 53 wherein the synthetic nucleotide sequence has a reduced number of one or more restriction
10 endonuclease recognition sites relative to the parent nucleic acid sequence.
62. A plasmid comprising the nucleic acid molecule of claim 53.
- 15 63. A plasmid which includes a sequence including SEQ ID NO:89, SEQ ID NO:90, or a sequence having at least 90% nucleic acid sequence identity thereto, or the complement thereof, which sequence encodes at least one selectable and/or screenable polypeptide.
- 20 64. The plasmid of claim 63 further comprising a multiple cloning region.
65. The plasmid of claim 63 further comprising another selectable or screenable polypeptide.
- 25 66. The plasmid of claim 63 or 65 wherein at least one selectable or screenable polypeptide comprises one or more protein destabilization sequences.
67. The plasmid of claim 63 wherein the sequence for the at least one
30 selectable and/or screenable polypeptide is not SEQ ID NO:41.
68. A synthetic nucleotide sequence of at least 100 nucleotides having a coding region for a selectable polypeptide which confers resistance to ampicillin, puromycin, hygromycin or neomycin, wherein the synthetic

nucleotide sequence has 90% or less nucleic acid sequence identity to a corresponding region of a parent nucleic acid sequence for the selectable polypeptide, wherein the decreased sequence identity is a result of different codons in the synthetic nucleotide sequence relative to the codons in the corresponding region in the parent nucleic acid sequence, wherein the synthetic nucleotide sequence has a reduced number of regulatory sequences relative to the average number of regulatory sequences resulting from random selections of codons at the sequences which differ between the synthetic nucleotide sequence and the parent nucleic acid sequence.

69. An isolated nucleic acid molecule encoding a selectable polypeptide and comprising a synthetic nucleotide sequence of at least 100 nucleotides having a coding region for the selectable polypeptide, wherein the synthetic nucleotide sequence has 90% or less nucleic acid sequence identity to a corresponding region in a parent nucleic acid sequence for the selectable polypeptide, wherein the decreased sequence identity is a result of different codons in the synthetic nucleotide sequence relative to the codons in the parent nucleic acid sequence, wherein the synthetic nucleotide sequence encodes a region of the selectable polypeptide with at least 85% amino acid sequence identity to the corresponding region of the selectable polypeptide encoded by the parent nucleic acid sequence, wherein the synthetic nucleotide sequence has a reduced number of regulatory sequences relative to the average number of regulatory sequences resulting from random selections of codons at the sequences which differ between the synthetic nucleotide sequence and the parent nucleic acid sequence, and wherein the isolated nucleic acid molecule, when expressed in a cell, confers resistance to ampicillin, puromycin, hygromycin or neomycin.

Figure 1

| <u>Amino Acid</u> | <u>Codon</u> |
|-------------------|------------------------------|
| Phe | UUU, UUC |
| Ser | UCU, UCC, UCA, UCG, AGU, AGC |
| Tyr | UAU, UAC |
| Cys | UGU, UGC |
| Leu | UUA, UUG, CUU, CUC, CUA, CUG |
| Trp | UGG |
| Pro | CCU, CCC, CCA, CCG |
| His | CAU, CAC |
| Arg | CGU, CGC, CGA, CGG, AGA, AGG |
| Gln | CAA, CAG |
| Ile | AUU, AUC, AUA |
| Thr | ACU, ACC, ACA, ACG |
| Asn | AAU, AAC |
| Lys | AAA, AAG |
| Met | AUG |
| Val | GUU, GUC, GUA, GUG |
| Ala | GCU, GCC, GCA, GCG |
| Asp | GAU, GAC |
| Gly | GGU, GGC, GGA, GGG |
| Glu | GAA, GAG |

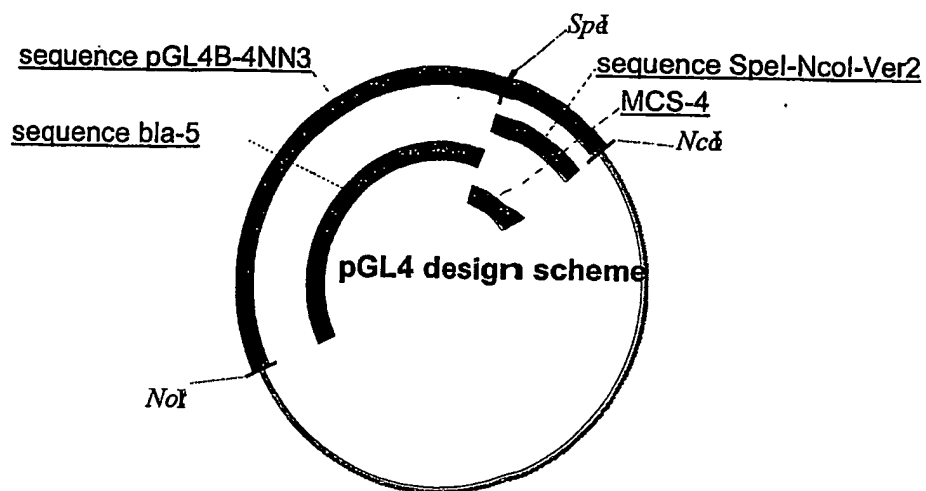


FIG. 2

SEQUENCE LISTING

<110> Promega Corporation
5 Wood, Keith
Wood, Monika
Almond, Brian
Paguio, Aileen
Fan, Frank

10

<120> Synthetic nucleic acid molecule and methods of preparation

<130> 341.034W01

15

<160> 97

<170> FastSEQ for Windows Version 4.0

20<210> 1

<211> 795

<212> DNA

<213> Unknown

25<220>

<223> Neo from neomycin gene from Promega's pCI-neo.

<400> 1

atgattgaac aagatggatt gcacgcaggt tctccggccg cttgggtgga gaggctattc 60
30ggctatgact gggcacaaca gacaatcggc tgctctgatg ccgccgtggt ccggctgtca 120
gcgcaggggc gcccggttct ttttgtcaag accgacctgt ccggtgccct gaatgaactg 180
caggacgagg cagcgcggct atcgtggctg gccacgacgg gcgttccttg cgcagctgtg 240
ctcgacgttg tcaactgaagc gggaagggac tggctgctat tgggcgaagt gccggggcag 300
gatctcctgt catctcacct tgctcctgcc gagaaagtat ccatcatggc tgatgcaatg 360
35cggcgggctgc atacgcttga tccggctacc tgccattcg accaccaagc gaaacatcgc 420
atcgagcgag cacgtactcg gatggaagcc ggtcttgtcg atcaggatga tctggacgaa 480
gagcatcagg ggctcgcgcc agccgaactg ttcgccaggc tcaaggcgcg catgcccgcac 540
ggcgaggatc tcgtcgtgac ccatggcgat gcctgcttgc cgaataatcat ggtggaaaat 600
ggccgctttt ctggattcat cgactgtggc cggctgggtg tggcggaccg ctatcaggac 660
40atagcgttgg ctaccggtga tattgtgaa gagcttggcg gcgaatgggc tgaccgcttc 720
ctcgtgcttt acggtatcgc cgctcccgat tcgcagcgca tcgccttcta tcgccttctt 780
gacgagttct tctga 795

2

<210> 2

<211> 264

<212> PRT

<213> Unknown

5

<220>

<223> Neo from neomycin gene from Promega's pCI-neo.

<400> 2

```

10Met Ile Glu Gln Asp Gly Leu His Ala Gly Ser Pro Ala Ala Trp Val
   1             5             10             15
   Glu Arg Leu Phe Gly Tyr Asp Trp Ala Gln Gln Thr Ile Gly Cys Ser
             20             25             30
   Asp Ala Ala Val Phe Arg Leu Ser Ala Gln Gly Arg Pro Val Leu Phe
15             35             40             45
   Val Lys Thr Asp Leu Ser Gly Ala Leu Asn Glu Leu Gln Asp Glu Ala
             50             55             60
   Ala Arg Leu Ser Trp Leu Ala Thr Thr Gly Val Pro Cys Ala Ala Val
   65             70             75             80
20Leu Asp Val Val Thr Glu Ala Gly Arg Asp Trp Leu Leu Leu Gly Glu
             85             90             95
   Val Pro Gly Gln Asp Leu Leu Ser Ser His Leu Ala Pro Ala Glu Lys
             100            105            110
   Val Ser Ile Met Ala Asp Ala Met Arg Arg Leu His Thr Leu Asp Pro
25             115            120            125
   Ala Thr Cys Pro Phe Asp His Gln Ala Lys His Arg Ile Glu Arg Ala
             130            135            140
   Arg Thr Arg Met Glu Ala Gly Leu Val Asp Gln Asp Asp Leu Asp Glu
   145            150            155            160
30Glu His Gln Gly Leu Ala Pro Ala Glu Leu Phe Ala Arg Leu Lys Ala
             165            170            175
   Arg Met Pro Asp Gly Glu Asp Leu Val Val Thr His Gly Asp Ala Cys
             180            185            190
   Leu Pro Asn Ile Met Val Glu Asn Gly Arg Phe Ser Gly Phe Ile Asp
35             195            200            205
   Cys Gly Arg Leu Gly Val Ala Asp Arg Tyr Gln Asp Ile Ala Leu Ala
             210            215            220
   Thr Arg Asp Ile Ala Glu Glu Leu Gly Gly Glu Trp Ala Asp Arg Phe
   225            230            235            240
40Leu Val Leu Tyr Gly Ile Ala Ala Pro Asp Ser Gln Arg Ile Ala Phe
             245            250            255
   Tyr Arg Leu Leu Asp Glu Phe Phe
             260

```

<210> 3

<211> 825

<212> DNA

5<213> Artificial Sequence

<220>

<223> A synthetic construct.

10<400> 3

```
ccactcagtg gccaccatga tcgagcagga cggcctgcac gccggcagcc ccgccgcctg 60
ggtaggagcgc ctgttcggct acgactgggc ccagcagacc atcggctgca gcgacgccgc 120
cgtgttcgcg ctgagcgccc agggccgccc cgtgctgttc gtgaagaccg acctgagcgg 180
cgccctgaac gagctgcagg acgaggccgc ccgcctgagc tggctggcca ccaccggcgt 240
15gccctgcgcc gccgtgctgg acgtggtgac cgaggccggc cgcgactggc tgctgctggg 300
cgaggtgccc ggccaggacc tgctgagcag ccacctggcc cccgccgaga aggtgagcat 360
catggccgac gccatgcgcc gcctgcacac cctggacccc gccacctgcc ccttcgacca 420
ccaggccaag caccgcatcg agcgcgccc caccgcatg gaggccggcc tggtaggacca 480
ggacgacctg gacgaggagc accagggcct ggcccccgcc gagctgttcg cccgcctgaa 540
20ggcccgcatg cccgacggcg aggacctggg ggtgaccac gccgacgct gcctgccc aa 600
catcatggtg gagaacggcc gcttcagcgg cttcatcgac tgcggccgcc tgggctggc 660
cgaccgctac caggacatcg ccctggccac ccgcgacatc gccgaggagc tgggctggc 720
gtgggcccgc cgcttcctgg tgctgtacgg catcgccgcc cccgacagcc agcgcacgc 780
cttctaccgc ctgctggacg agttcttcta ataaccagtc tctgg 825
```

25

<210> 4

<211> 825

<212> DNA

<213> Artificial Sequence

30

<220>

<223> A synthetic construct.

<400> 4

```
35ccactcagtg gccaccatga tcgagcagga cggcctgcac gccggcagcc ccgccgcctg 60
ggtaggagcgc ctgttcggct acgactgggc ccagcagacc atcggctgca gcgacgccgc 120
cgtgttcgcg ctgagcgccc agggccgccc cgtgctgttc gtgaagaccg acctgagcgg 180
cgccctgaac gagctgcagg acgaggccgc ccgcctgagc tggctggcca ccaccggcgt 240
ggccctgcgcc gccgtgctgg acgtggtgac cgaggccggc cgcgactggc tgctgctggg 300
40cgaggtgccc ggccaggacc tgctgagcag ccacctggcc cccgccgaga aggtgagcat 360
catggccgac gccatgcgcc gcctgcacac cctggacccc gccacctgcc ccttcgacca 420
ccaggccaag caccgcatcg agcgcgccc caccgcatg gaggccggcc tggtaggacca 480
aaacacacctg aacaaqaaqc accaaqaccc aacccccacc aaactattca cccacctaaa 540
```

ggccccgatg cccgacggcg aggacctggt ggtgaccac ggcgacgcct gcctgccc aa 600
catcatggtg gagaacggcc gcttcagcgg cttcatcgac tgcggccgcc tgggcgtggc 660
cgaccgctac caggacatcg ccctggccac ccgcgacatc gccgaggagc tgggcggcga 720
gtgggcccgc cgcttcctgg tgctgtacgg catcgccgcc cccgacagcc agcgcatcgc 780
5cttctaccgc ctgctggacg agttcttcta ataaccagtc tctgg 825

<210> 5

<211> 818

<212> DNA

10<213> Artificial Sequence

<220>

<223> A synthetic construct.

15<400> 5

cctgcaggcc accatgatcg aacaagacgg cctccatgct ggcagtcccg cagcttgggt 60
cgaacgcttg ttcgggtacg actgggccc a gcagaccatc ggatgtagcg atgcggccgt 120
gttcgcgtcta agcgctcaag gccggcccgt gctgttcgtg aagaccgacc tgagcggcgc 180
cctgaacgag cttcaagacg aggctgccc cctgagctgg ctggccacca ccggtgtacc 240
20ctgcgcgct gtgttgatg ttgtgaccga agcggcccg gactggctgc tgctgggcga 300
ggtcacctggc caggatctgc tgagcagcca ccttgcccc gctgagaagg tttccatcat 360
ggccgatgca atgcggcgcc tgcacaccct ggaccccgct acatgcccct tcgaccacca 420
ggctaagcat cggatcgagc gtgctcggac ccgcatggag gccggcctgg tggaccagga 480
cgacctggac gaggagcatc agggcctggc ccccgctgaa ctgttcgccc gcctgaaagc 540
25ccgcatgccg gacggtgagg acctggttgt gacacatggt gatgcctgcc tccctaacat 600
catggctcag aatggccgct tctccggctt catcgactgc ggtcgcctag gagttgccga 660
ccgctaccag gacatcgccc tggccacccg cgacatcgct gaggagcttg gcggcgagtg 720
ggccgaccgc ttcttagtct tgtacggcat cgcagctccc gacagccagc gcatcgctt 780
ctaccgcctg ctcgacgagt tcttttaatg agcttaag 818

30

<210> 6

<211> 1024

<212> DNA

<213> Escherichia coli

35

<400> 6

atgaaaaagc ctgaactcac cgcgacgtct gtcgagaagt ttctgatcga aaagttcgac 60
agcgtctccg acctgatgca gctctcggag gggaagaat ctcggtgctt cagcttcgat 120
gtaggagggc gtggatatgt cctgcgggta aatagctgcg ccgatggttt ctacaaagat 180
40cgttatgtt atcggaactt tgcacggcc gcgctcccga ttccggaagt gcttgacatt 240
ggggaattca gcgagagcct gacctattgc atctcccgc gtgcacaggg tgtcacgttg 300
caagacctgc ctgaaaccga actgcccgtt gttctgcagc cggtcgcgga ggccatggat 360
accatcacta cagccatctt taaccaaacg aacgaattca acccattcag acccaaaqa 420

atcgggtcaat acactacatg gcgtgatttc atatgcgcga ttgctgatcc ccatgtgtat 480
 cactggcaaa ctgtgatgga cgacaccgtc agtgcgtccg tcgcgcaggc tctcgatgag 540
 ctgatgcttt gggccgagga ctgccccgaa gtccggcacc tcgtgcacgc ggatttcggc 600
 tccaacaatg tcctgacgga caatggccgc ataacagcgg tcattgactg gagcgaggcg 660
 5atgttcgggg attcccaata cgaggtcgcc aacatcttct tctggaggcc gtggttggt 720
 tgtatggagc agcagacgcg ctacttcgag cggaggcatc cggagcttgc aggatcgccg 780
 cggtccggg cgtatatgct ccgcattggt cttgaccaac tctatcagag cttgggtgac 840
 ggcaatttcg atgatgcagc ttgggcgcag ggtcgatgcg acgcaatcgt ccgatccgga 900
 gccgggactg tcgggcgtac acaaatcgcc cgcagaagcg cggccgtctg gaccgatggc 960
 10tgtgtagaag tactcgccga tagtggaac cgacgcccc gactcgtcc gagggcaaag 1020
 gaat 1024

<210> 7

<211> 341

15<212> PRT

<213> Escherichia coli

<400> 7

Met Lys Lys Pro Glu Leu Thr Ala Thr Ser Val Glu Lys Phe Leu Ile
 20 1 5 10 15
 Glu Lys Phe Asp Ser Val Ser Asp Leu Met Gln Leu Ser Glu Gly Glu
 20 25 30
 Glu Ser Arg Ala Phe Ser Phe Asp Val Gly Gly Arg Gly Tyr Val Leu
 35 40 45
 25Arg Val Asn Ser Cys Ala Asp Gly Phe Tyr Lys Asp Arg Tyr Val Tyr
 50 55 60
 Arg His Phe Ala Ser Ala Ala Leu Pro Ile Pro Glu Val Leu Asp Ile
 65 70 75 80
 Gly Glu Phe Ser Glu Ser Leu Thr Tyr Cys Ile Ser Arg Arg Ala Gln
 30 85 90 95
 Gly Val Thr Leu Gln Asp Leu Pro Glu Thr Glu Leu Pro Ala Val Leu
 100 105 110
 Gln Pro Val Ala Glu Ala Met Asp Ala Ile Ala Ala Ala Asp Leu Ser
 115 120 125
 35Gln Thr Ser Gly Phe Gly Pro Phe Gly Pro Gln Gly Ile Gly Gln Tyr
 130 135 140
 Thr Thr Trp Arg Asp Phe Ile Cys Ala Ile Ala Asp Pro His Val Tyr
 145 150 155 160
 His Trp Gln Thr Val Met Asp Asp Thr Val Ser Ala Ser Val Ala Gln
 40 165 170 175
 Ala Leu Asp Glu Leu Met Leu Trp Ala Glu Asp Cys Pro Glu Val Arg
 180 185 190

His Leu Val His Ala Asp Phe Gly Ser Asn Asn Val Leu Thr Asp Asn
 195 200 205
 Gly Arg Ile Thr Ala Val Ile Asp Trp Ser Glu Ala Met Phe Gly Asp
 210 215 220
 5Ser Gln Tyr Glu Val Ala Asn Ile Phe Phe Trp Arg Pro Trp Leu Ala
 225 230 235 240
 Cys Met Glu Gln Gln Thr Arg Tyr Phe Glu Arg Arg His Pro Glu Leu
 245 250 255
 Ala Gly Ser Pro Arg Leu Arg Ala Tyr Met Leu Arg Ile Gly Leu Asp
 10 260 265 270
 Gln Leu Tyr Gln Ser Leu Val Asp Gly Asn Phe Asp Asp Ala Ala Trp
 275 280 285
 Ala Gln Gly Arg Cys Asp Ala Ile Val Arg Ser Gly Ala Gly Thr Val
 290 295 300
 15Gly Arg Thr Gln Ile Ala Arg Arg Ser Ala Ala Val Trp Thr Asp Gly
 305 310 315 320
 Cys Val Glu Val Leu Ala Asp Ser Gly Asn Arg Arg Pro Ser Thr Arg
 325 330 335
 Pro Arg Ala Lys Glu
 20 340

<210> 8

<211> 1056

<212> DNA

25<213> Artificial Sequence

<220>

<223> A synthetic construct.

30<400> 8

ccactcagtg gccaccatga agaagcccga gctgaccgcc accagcgtgg agaagttcct 60
 gatcgagaag ttcgacagcg tgagcgacct gatgcagctg agcgagggcg aggagagccg 120
 cgccttcagc ttcgacgtgg gcggccgcgg ctacgtgctg cgcgtgaaca gctgcgccga 180
 cggcttctac aaggaccgct acgtgtaccg ccacttcgcc agcgccgccc tgcccatccc 240
 35cgaggtgctg gacatcggcg agttcagcga gagcctgacc tactgcatca gccgccgcgc 300
 ccagggcgctg accctgcagg acctgcccga gaccgagctg cccgccgtgc tgcagcccgt 360
 ggccgaggcc atggacgcca tcgccgccgc gcacctgagc cagaccagcg gcttcggccc 420
 cttcggcccc cagggcatcg gccagtacac cacctggcgc gacttcatct gcgccatcgc 480
 cgacccccac gtgtaccact ggcagaccgt gatggacgac accgtgagcg ccagcgtggc 540
 40ccagggcctg gacgagctga tgctgtgggc cgaggactgc cccgaggtgc gccacctggt 600
 gcacgccgac ttggcgagca acaacgtgct gaccgacaac ggccgcatca ccgccgtgat 660
 cgactggagc gaggccatgt tcggcgacag ccagtacgag gtggccaaca tcttcttctg 720
 gcgcccttgg ctggcctgca tggagcagca gaccgcctac ttcgagcgc qccacccccq 780

7

gctggccggc agcccccgcc tgcgcgccta catgctgcgc atcggcctgg accagctgta 840
ccagagcctg gtggacggca acttcgacga cgccgcctgg gccagggcc gctgcgacgc 900
catcgtgcgc agcggcgccg gcaccgtggg ccgcacccag atcgcccgcc gcagcgccgc 960
cgtgtggacc gacggctgcg tggaggtgct ggccgacagc ggcaaccgcc gcccagcac 1020
5ccgccccgc gccaaaggagt aataaccagc tcttgg 1056

<210> 9

<211> 1056

<212> DNA

10<213> Artificial Sequence

<220>

<223> A synthetic construct.

15<400> 9

ccactccgtg gccaccatga agaagcccga gctgaccgct accagcggtg aaaaatttct 60
catcgagaag ttcgacagtg tgagcgacct gatgcagttg tcggagggcg aagagagccg 120
agccttcagc ttcgatgtcg gcggacgcgg ctatgtactg cgggtgaata gctgcgctga 180
tggcttctac aaagaccgct acgtgtaccg ccacttcgcc agcgtgcac taccatccc 240
20cgaagtgttg gacatcgcg agttcagcga gaggctgaca tactgcatca gtagacgcgc 300
ccaaggcggtt actctccaag acctccccga aacagagctg cctgctgtgt tacagcctgt 360
cgccgaagct atggatgcta ttgcccgccg cgacctcagt caaaccagcg gcttcggccc 420
attcggggcc caaggcatcg gccagtacac aacctggcgg gatttcattt gcgccattgc 480
tgatccccat gtctaccact ggcagaccgt gatggacgac accgtgtccg ccagcgtagc 540
25tcaagccctg gacgaactga tgctgtgggc cgaagactgt cccgaggtgc gccacctcgt 600
ccatgccgac ttcggcagca acaacgtcct gaccgacaac ggccgcatca ccgccgtaat 660
cgactggtcc gaagctatgt tcggggacag tcagtacgag gtggccaaca tcttcttctg 720
gcggccctgg ctggcttgca tggagcagca gactcgctac ttcgagcgcc ggcatcccga 780
gctggccggc agccctcgtc tgcgagccta catgctgcgc atcggcctgg atcagctcta 840
30ccagagcctc gtggacggca acttcgacga tgctgcctgg gctcaaggcc gctgcgatgc 900
catcgtccgc agcggggccg gcaccgtcgg tcgcacacaa atcgctcgcc ggagcgccgc 960
cgtatggacc gacggctgcg tcgaggtgct ggccgacagc ggcaaccgcc ggcccagtac 1020
acgaccgcgc gctaaggagt agtaaccagg ctctgg 1056

35<210> 10

<211> 1048

<212> DNA

<213> Artificial Sequence

40<220>

<223> A synthetic construct.

<400> 10
cctgcaggcc accatgaaga agcccgagct gaccgctacc agcgttgaaa aattttctcat 60
cgagaagttc gacagtgtga gcgacctgat gcagttgtcg gagggcgaag agagccgagc 120
cttcagcttc gatgtcggcg gacgcggcta tgtactgagg gtgaatagct gcgctgatgg 180
5ctttctacaaa gaccgctacg tgtaccgcca cttcgccagc gctgcactac ccatccccga 240
agtgttggac atcggcgagt tcagcgagag cctgacatac tgcatacagta gacgcgcccc 300
aggcgttact ctccaagacc tccccgaaac agagctgcct gctgtgttac agcctgtcgc 360
cgaagctatg gatgctattg ccgccgccga cctcagtcaa accagcggct tcggcccatt 420
cgggccccc aa ggcatcggcc agtacacaac ctggcgggat ttcatttgcg ccattgctga 480
10tccccatgtc taccactggc agaccgtgat ggacgacacc gtgtccgcca gcgtagctca 540
agccctggac gaactgatgc tgtggggccga agactgtccc gaggtgcgcc acctcgacca 600
tgccgacttc ggagcaaca acgtcctgac cgacaacggc cgcatacccg ccgtaatcga 660
ctggctccga gctatgttcg gggacagtca gtacgaggtg gccaacatct tcttctggcg 720
gccctggctg gcttgcattg agcagcagac tcgctacttc gagcgccggc atcccagagct 780
15ggccggcagc cctcgtctgc gagcctacat gctgcgcata ggccctggatc agctctacca 840
gagcctcgtg gacggcaact tcgacgatgc tgcctgggct caaggccgct gcgatgccat 900
cgtccgcagc ggggccggca ccgtcggctc caccacaaatc gctcgccgga gcgcgcgct 960
atggaccgac ggctgcgtcg aggtgctggc cgacagcggc aaccgcggc ccagtagacg 1020
accgcgcgt aaggagtagt aacttaag 1048

20
<210> 11
<211> 1174
<212> DNA
<213> Artificial Sequence

25
<220>
<223> A synthetic construct.

<400> 11
3Oggatccgttt gcgtattggg cgtcttccg ctgatctgag cagcaccatg gcctgaaata 60
acctctgaaa gaggaacttg gttagctacc ttctgaggcg gaaagaacca gctgtggaat 120
gtgtgtcagt tagggtgtgg aaagtcccca ggctcccag caggcagaag tatgcaaagc 180
atgcatctca attagtcagc aaccaggtgt ggaaagtccc caggctcccc agcaggcaga 240
agtatgcaaa gcatgcatct caattagtca gcaaccatag tcccgcccct aactccgccc 300
35atcccccccc taactccgcc cagttccgcc cattctccgc cccatggctg actaattttt 360
tttatttatg cagaggccga ggccgctct gcctctgagc tattccagaa gtagtgagga 420
ggcttttttg gaggcctagg cttttgcaaa aagctcgatt cttctgacac tagcgccacc 480
atgaccgagt acaagcctac cgtgcgcctg gccactcgcg atgatgtgcc ccgcgcgctc 540
cgcaactctgg ccgccgcttt cgccgactac ccgctaccc ggacacaccg ggaccccgac 600
40cggcacatcg agcgtgtgac agagttgcag gagctgttcc tgaccgcgct cgggctggac 660
atcggaagcgt ggtgggtagc cgacgacggc gggccgctgg ccgtgtggac tcccccgag 720
agcgttgagg ccggcgccgt gttcgccgag atcgcccccc gaatggccga gctgagcggc 780
agccgcctgg ccgccagca gcaaatggag ggccgtgctt cccccatcg tcccaaggag 840

```

cctgcctgggt ttctggccac tgtaggagtg agccccgacc accagggcaa gggcttgggc 900
agcgccgctcg tgttgcccgg cgtagaggcc gccgaacgcg ccggtgtgcc cgcctttctc 960
gaaacaagcg caccaagaaa ccttccattc tacgagcgcc tgggcttcac cgtgaccgcc 1020
gatgtcgagg tgcccgaggg acctaggacc tgggtgatga cacgaaaacc tggcgccata 1080
5tgatctagaa ccggtcatgg ccgcaataaa atatctttat tttcattaca tctgtgtgtt 1140
ggttttttgt gtgttcgaac tagatgctgt cgac 1174

```

<210> 12

<211> 1776

10<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

15

<400> 12

```

atggcttcca aggtgtacga ccccgagcaa cgcaaacgca tgatcactgg gcctcagtgg 60
tgggctcgtc gcaagcaaat gaacgtgctg gactccttca tcaactacta tgattccgag 120
aagcacgccc agaacgcggt gatttttctg catggtaacg ctgcctccag ctacctgtgg 180
20agggcacgtc tgccctacat cgagcccgtg gctagatgca tcatccctga tctgatcgga 240
atgggtaagt ccggcaagag cgggaatggc tcatatcgcc tcctggatca ctacaagtac 300
ctcacgctt ggttcgagct gctgaacctt ccaaagaaaa tcatctttgt gggccacgac 360
tgggggggct gtctggcctt tcactactcc tacgagcacc aagacaagat caaggccatc 420
gtccatgctg agagtgtcgt ggacgtgatc gagtcctggg acgagtggcc tgacatcgag 480
25gaggatatcg ccctgatcaa gagcgaagag ggcgagaaaa tgggtgcttga gaataacttc 540
ttcgtcgaga ccatgctccc aagcaagatc atgcggaaac tggagcctga ggagttcgct 600
gcctacctgg agccattcaa ggagaagggc gaggttagac ggccctacct ctctggcct 660
cgcgagatcc ctctcgtaa gggaggcaag cccgacgtcg tccagattgt ccgcaactac 720
aacgcctacc ttcgggccag cgacgatctg cctaagatgt tcatcgagtc cgaccctggg 780
30ttcttttcca acgctattgt cgagggagct aagaagtcc ctaacaccga gttcgtgaag 840
gtgaagggcc tccacttcag ccaggaggac gctccagatg aaatgggtaa gtacatcaag 900
agcttcgtgg agcgctgct gaagaacgag cagaccggtg gtgggagcgg aggtggcgga 960
tcagggtggc gaggtccgg agggattgaa caagatggat tgcacgcagg ttctccggcc 1020
gcttgggtgg agaggctatt cggctatgac tgggcacaac agacaatcgg ctgctctgat 1080
35gccgccgtgt tccggtgtc agcgagggg cgcccggtt tttttgtcaa gaccgacctg 1140
tccggtgccc tgaatgaact gcaggacgag gcagcgcggc tatcgtggct ggccacgacg 1200
ggcgttcctt gcgcagctgt gctcgacgtt gtcactgaag cgggaaggga ctggctgcta 1260
ttgggcgaag tgccggggca ggatctcctg tcatctcacc ttgctcctgc cgagaaagta 1320
tccatcatgg ctgatgcaat gcggcggtg catacgcttg atccggctac ctgccattc 1380
40gaccaccaag cgaaacatcg catcgagcga gcacgtactc ggatggaagc cggctctgtc 1440
gatcaggatg atctggacga agagcatcag gggctcgcg cagccgaact gttcgccagg 1500
ctcaaggcgc gcatgcccga cggcgaggat ctctctgtga cccatggcga tgccctgctt 1560
ccgaatatca tgggtggaaa tggccgctt tctggattca tcgactgtgg ccggctgggt 1620

```

10

```

gtggcggacc gctatcagga catagcgttg gctaccctg atattgctga agagcttggc 1680
ggcgaatggg ctgaccgctt cctcgtgctt tacggatcg ccgctcccga ttcgcagcgc 1740
atcgcttct atcgcttct tgacgagttc ttctaa 1776

```

5<210> 13

<211> 1776

<212> DNA

<213> Artificial Sequence

10<220>

<223> A synthetic construct.

<400> 13

```

atgattgaac aagatggatt gcacgcaggt tctccggccg cttgggtgga gaggctattc 60
15ggctatgact gggcacaaca gacaatcggc tgctctgatg ccgccgtgtt ccggctgtca 120
gcgcaggggc gcccggttct ttttgtcaag accgacctgt ccggtgccct gaatgaactg 180
caggacgagg cagcgcggct atcgtggctg gccacgacgg gcgttccttg cgcagctgtg 240
ctcgacgttg tactgaagc gggaaaggac tggctgctat tgggcgaagt gccggggcag 300
gatctcctgt catctcacct tgctcctgcc gagaaagtat ccatcatggc tgatgcaatg 360
20cggcggctgc atacgcttga tccggctacc tgccattcg accaccaagc gaaacatcgc 420
atcgagcgag cagtgactcg gatggaagcc ggtcttgctg atcaggatga tctggacgaa 480
gagcatcagg ggtcgcgcc agccgaactg ttcgccaggc tcaaggcgcg catgcccagc 540
ggcgaggatc tcgtcgtgac ccatggcgat gcctgcttgc cgaatatcat ggtggaaaat 600
ggcgcgtttt ctggattcat cgactgtggc cggctgggtg tggcggaccg ctatcaggac 660
25atagcgttg ctaccctgta tattgctgaa gagcttggcg gcgaatgggc tgaccgcttc 720
ctcgtgcttt acggtatcgc cgtcccgat tcgcagcgca tcgccttcta tcgccttctt 780
gacgagttct tcaccgggtg tgggagcggg ggtggcggat cagggtggcg aggtccgga 840
ggggcttcca aggtgtacga ccccgagcaa cgaacacgca tgatcactgg gcctcagtgg 900
tgggctcgtc gcaagcaaat gaacgtgctg gactccttca tcaactacta tgattccgag 960
30aagcacgccc agaacgccgt gatttttctg catggtaacg ctgcctccag ctacctgtgg 1020
aggcacgtcg tgcctcacat cgagcccgtg gctagatgca tcatccctga tctgatcgga 1080
atgggtaagt ccggcaagag cgggaatggc tcatatcgcc tcctggatca ctacaagtac 1140
ctcacgcgtt ggttcgagct gctgaacctt ccaaagaaaa tcatctttgt gggccacgac 1200
tggggggcct gtctggcctt tcactactcc tacgagcacc aagacaagat caaggccatc 1260
35gtccatgctg agagtgtcgt ggacgtgatc gagtcctggg acgagtggcc tgacatcgag 1320
gaggatatcg ccctgatcaa gagcgaagag ggcgagaaaa tgggtgcttga gaataacttc 1380
ttcgtcgaga ccatgctccc aagcaagatc atgcggaaac tggagcctga ggagtctcgt 1440
gcctacctgg agccattcaa ggagaagggc gaggttagac ggccacctc ctctggcct 1500
cgcgagatcc ctctcgttaa gggaggcaag cccgacgtcg tccagattgt ccgcaactac 1560
40aacgcctacc ttcgggccag cgacgatctg cctaagatgt tcatcgagtc cgaccctggg 1620
ttcttttcca acgtattgt cgaggagct aagaagttcc ctaacaccga gttcgtgaag 1680
gtgaaggggc tccacttcag ccaggaggac gctccagatg aaatgggtaa gtacatcaag 1740
agcttcgtgg agcgcgtgct gaagaacgag cagtaa 1776

```

<210> 14

<211> 1653

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 14

```
10atggccgatg ctaagaacat taagaagggc cctgctccct tctaccctct ggaggatggc 60
   accgctggcg agcagctgca caaggccatg aagaggatat ccctggtgcc tggcaccatt 120
   gccttcaccg atgccacat tgaggtggac atcacctatg ccgagtactt cgagatgtct 180
   gtgcgccctgg ccgaggccat gaagaggtag gccctgaaca ccaaccaccg catcgtggtg 240
   tgctctgaga actctctgca gttcttcata ccagtgtctg gcgccctgtt catcggagtg 300
15gccgtggccc ctgctaacga catttacaac gagcgcgagc tgctgaacag catgggcatt 360
   tctcagccta ccgtggtgtt cgtgtctaag aagggcctgc agaagatcct gaacgtgcag 420
   aagaagctgc ctatcatcca gaagatcatc atcatggact ctaagaccga ctaccagggc 480
   ttccagagca tgtacacatt cgtgacatct catctgcctc ctggcctcaa cgagtacgac 540
   ttctgtgccag agtctttcga cagggacaaa accattgccc tgatcatgaa cagctctggg 600
20tctaccggcc tgcctaaggg cgtggccctg cctcatcgca ccgcctgtgt gcgcttctct 660
   cacgcccgcg accctatttt cggcaaccag atcatccccg acaccgctat tctgagcgtg 720
   gtgccattcc accacggcct cggcatgttc accaccctgg gctacctgat ttggggcttt 780
   cgggtggtgc tgatgtaccg cttcgaggag gagctgttcc tgcgcagcct gcaagactac 840
   aaaattcagt ctgccctgct ggtgccaacc ctgttcagct tcttcgctaa gagcaccctg 900
25atcgacaagt acgacctgtc taacctgcac gagattgcct ctggcgggcg cccactgtct 960
   aaggagggtg gcgaagccgt ggccaagcgc tttcatctgc caggcatccg ccagggctac 1020
   ggcctgaccg agacaaccag cgccattctg attacccag agggcgacga caagcctggc 1080
   gccgtgggca aggtggtgcc attcttcgag gccaaagggtg tggacctgga caccggcaag 1140
   accctgggag tgaaccagcg cggcgagctg tgtgtgcgcg gccctatgat tatgtccggc 1200
30tacgtgaata accctgaggg cacaaacgcc ctgatcgaca aggacggctg gctgcactct 1260
   ggcgacattg cctactggga cgaggacgag cacttcttca tcgtggaccg cctgaagtct 1320
   ctgatcaagt acaagggtcta ccaggtggcc ccagccgagc tggagtctat cctgctgcag 1380
   caccctaaca ttttcgacgc cggagtggcc ggcctgcccg acgacgatgc cggcgagctg 1440
   cctgccgcg tcgtcgtgct ggaacacggc aagaccatga ccgagaagga gatcgtggac 1500
35tatgtggcca gccaggtgac aaccgcc aag aagctgcgcg gcggagtggg gtctgtggac 1560
   gaggtgcca agggcctgac cggcaagctg gacgcccga agatccgca gatcctgac 1620
   aaggctaaga aaggcggcaa gatcgccgtg taa 1653
```

<210> 15

40<211> 597

<212> DNA

<213> Streptomyces sp.

<400> 15
 atgaccgagt acaagcccac ggtgcgcctc gccacccgcg acgacgtccc ccggggccgta 60
 cgcaccctcg ccgccgcgtt cgcgcactac ccgcgccacgc gccacaccgt cgacccggac 120
 cgccacatcg agcgggtcac cgagctgcaa gaactcttcc tcacgcgcgt ccgggtctcgac 180
 5atcggcaagg tgtgggtcgc ggacgacggc gccgcggtgg cgggtctggac cacgccggag 240
 agcgtcgaag ccgggggcgtt gttcgcgcgag atcggcccgc gcatggccga gttgagcggg 300
 tcccggctgg ccgcgcagca acagatggaa ggcctcctgg ccgcgcaccg gcccaaggag 360
 ccgcgctggg tcctggccac cgtcggcgtg tcgcccgaac accagggcaa ggggtctgggc 420
 agcgcgcgtc tgctccccgg agtggaggcg gccgagcgcg ccgggggtgcc cgccttcctg 480
 10gagacctccg cgcgccgcaa cctccccctt tacgagcggc tcggcttcac cgtcaccgcc 540
 gacgtcgagg tgcccgaagg accgcgcacc tgggtcatga cccgcaagcc cgggtgcc 597

<210> 16
 <211> 1672
 15<212> DNA
 <213> Artificial Sequence

<220>
 <223> A synthetic construct.

20
 <400> 16
 aaagccacca tggaggacgc caagaacatc aagaaggggc ccgccccctt ctacccccctg 60
 gaggacggca ccgccggcga gcagctgcac aaggccatga agcgtctacgc cctgggtgcc 120
 ggcaccatcg ccttcaccga cgcacacatc gaggtggaca tcacctacgc cgagtacttc 180
 25gagatgagcg tgcgcctggc cgaggccatg aagcgtctacg gcctgaacac caaccaccgc 240
 atcgtgggtg gcagcgagaa cagcctgcag ttcttcatgc ccgtgctggg cgcctgttgc 300
 atcggcgtgg ccgtggcccc cgccaacgac atctacaacg agcgcgagct gctgaacagc 360
 atgggcatca gccagcccac cgtggtgttc gtgagcaaga agggcctgca gaagatcctg 420
 aacgtgcaga agaagctgcc catcatccag aagatcatca tcattggacag caagaccgac 480
 30taccagggct tccagagcat gtacaccttc gtgaccagcc acctgcccc cggcttcaac 540
 gagtacgact tcgtgcccga gagcttcgac cgcgacaaga ccatcgccct gatcatgaac 600
 agcagcggca gcaccggcct gcccgaaggc gtggccctgc cccaccgcac cgctgctgtg 660
 cgcttcagcc acgcccgcga ccccatcttc ggcaaccaga tcatccccga caccgccatc 720
 ctgagcgtgg tgcccttcca ccacggcttc ggcattgtca ccacctggg ctacctgac 780
 35tgcggttcc gcgtggtgct gatgtaccgc ttcgaggagg agctgttcct gcgcagcctg 840
 caggactaca agatccagag cgccctgctg gtgcccaccc tgttcagctt cttcgccaag 900
 agcaccctga tcgacaagta cgacctgagc aacctgcacg agatcgccag cggcgggcgc 960
 cccctgagca aggagggtgg cgaggccgtg gccaaagcgt tcacacctgc cggcatccgc 1020
 cagggtctac gcctgaccga gaccaccagc gccatcctga tcacccccga gggcgacgac 1080
 40aagcccggcg ccgtgggcaa ggtggtgccc ttcttcgagg ccaagggtgg ggacctggac 1140
 accggcaaga ccctgggcgt gaaccagcgc ggcgagctgt gcgtgcgcgg ccccatgac 1200
 atgagcggct acgtgaacaa ccccgaggcc accaacgccc tgatcgacaa ggacggctgg 1260
 ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgc 1320

13

ctgaagagcc tgatcaagta caagggctac caggtggccc ccgccgagct ggaagagcatc 1380
 ctgctgcagc accccaacat cttcgacgcc ggcgtggccg gcctgcccga cgaagacgcc 1440
 ggcgagctgc ccgccgccgt ggtggtgctg gagcacggca agaccatgac cgaagaaggag 1500
 atcgtggact acgtggccag ccaggtgacc accgccaaga agctgcgcgg cggcgtggtg 1560
 5ttcgtggacg aggtgcccga ggcctgacc ggcaagctgg acgcccga ga tccgcgag 1620
 atcctgatca aggccaagaa gggcggcaag atcgccgtgt aataattcta ga 1672

<210> 17

<211> 1672

10<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

15

<400> 17

aaagccacca tggaggagcg caagaacatc aagaagggcc cagcgccatt ctacccccctg 60
 gaggacggca ccgccggcga gcagctgcac aaggccatga agcgctacgc cctggtgccc 120
 ggcaccatcg ccttcacga cgcacatata gaggtggaca tcacctacgc cgaagtacttc 180
 20gagatgagcg ttcggctggc agaggctatg aagcgctatg ggctgaacac caaccatcgc 240
 atcgtggtgt gcagcgagaa cagcttgacg ttcttcatgc ccgtgttggg tgcctgttgc 300
 atcggcgtgg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360
 atgggcatca gccagccac cgtcgtatc gtgagcaaga aagggtgca aaagatcctg 420
 aacgtgcaaa agaagctgcc catcatcaa aagatcatca tcatggacag caagaccgac 480
 25taccagggtc tccaaagcat gtacacctc gtgaccagcc atttgccgc cggcttcaac 540
 gagtacgact tcgtgcccga gagcttcgac cgcgacaaga ccatcgccct gatcatgaac 600
 agtagtggca gtaccggctt acctaagggc gtggccctac cgcaccgcac cgcctgtgtc 660
 cgattcagtc atgcccgcga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
 ctgagcgtgg tgccatttca ccacggcttc ggcattgtca ccaccctggg ctacttgatc 780
 30tcggcgttcc ggtcgtgct gatgtaccgc ttcgaggagg agctattctt cgcgagcttg 840
 caagactaca agattcaaag cgccctgctg gtgcccaccc tgttcagttt cttcgccaag 900
 agcaccctga tcgacaagta cgacctgacg aacctgcacg agatcgccag cggcgccgcc 960
 ccgctcagca aggaggtggg cgaggccgtg gccaaagcgt tccacctgcc aggcacccgc 1020
 cagggtacg gcctgaccga gacaaccagc gccattctga tcacccccga gggggacgac 1080
 35aagcctggcg cagtaggcaa ggtggtgccc ttcttcgagg ctaaggtggt ggaacctggac 1140
 accggtaaaa ccctgggtgt gaaccagcgc ggcgagctgt gcgtccgtgg ccccatgatc 1200
 atgagcggct acgttaacaa ccccgaggct acaaacgccc tgatcgacaa ggaaggctgg 1260
 ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
 ctgaagagcc tgatcaata caagggctac caggtagccc cagccgaact ggaagagcatc 1380
 40ctgctgcagc accccaacat cttcgacgcc ggggtcgccg gcctgcccga cgaagatgcc 1440
 ggcgagctgc ccgcccgagt cgtggtgctg gagcacggta aaaccatgac cgaagaaggag 1500
 atcgtggact atgtggccag ccaggttaca accgccaaga agctgcgcgg cggcgtggtg 1560
 ttcgtggacg aggtgcctaa aggcctgacg ggcaagttgg acgcccga ga tccgcgag 1620

14

attctgatca aggccaagaa gggcggcaag atcgccgtgt aataattcta ga 1672

<210> 18

<211> 1672

5<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

10

<400> 18

aaagccacca tggaagatgc caaaaacatt aagaagggcc cagcgccatt ctacccactg 60
gaggacggca ccgccggcga gcagctgcac aaagccatga agcgctacgc cctggtgccc 120
ggcaccatcg cctttaccga cgcacatcgc gaggtggaca tcacctacgc cgagtacttc 180
15gagatgagcg ttccggctggc agaggctatg aagcgctatg ggctgaatac caaccatcgc 240
atcggtggtg gcagcgagaa tagcttgagc ttcttcacgc ccgtggtggg tgccctgttc 300
atcggtgtgg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360
atgggcatca gccagcccac cgtcgtattc gtgagcaaga aagggtgca aaagatcctc 420
aacgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
20taccagggct tccaaagcat gtacaccttc gtgaccagcc atttgccacc cggtttcaac 540
gagtacgact tcgtgcccga gagcttcgac cgggacaaaa ccatcgccct gatcatgaac 600
agtagtggca gtaccggatt gccaagggc gtagccctac cgcaccgcac cgctgtgttc 660
cgattcagtc atgcccgcga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
ctcagcgtgg tgccatttca ccacggcttc ggcattgtca ccacgctggg ctacttgatc 780
25tcgcgctttc gggctcgtgt catgtaccgc ttccaggagg agctattctt gcgcagcttg 840
caagactata agattcaaag cgccctgctg gtgcccacac tggtcagctt cttcgccaaag 900
agcactctca tcgacaagta cgacctgagc aacctgcacg agatcgccag cggcggggcg 960
ccgtcagca aggaggtggg cgaggccgtg gccaaagcgt tccacctacc aggcattccg 1020
cagggtacg gcctgacaga aacaaccagc gccattctga tcacccccga aggggacgac 1080
30aagcctggcg cagtaggcaa ggtgggtgccc ttcttcgagg ctaagggtgt ggacttggac 1140
accggtgaaga ccctgggtgt gaaccagcgc ggcgagctgt gcgtccgtgg ccccatgatc 1200
atgagcgggt acgttaacaa ccccgaggt acaaacgctc tcatcgacaa ggacggctgg 1260
ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
ctgaagagcc tgatcaaata caagggtac caggtagccc cagccgaact ggagagcatc 1380
35ctgctgcaac accccaacat cttcgacgcc ggggtcgccg gcctgcccga cgacgatgcc 1440
ggcgagctgc ccgccgcagt cgtcgtgctg gagcacggta aaacctgac cgagaaggag 1500
atcggtgact atgtggccag ccagggtaca accgccaaga agctgcgcgg tgggtgtgtg 1560
ttcgtggacg aggtgcctaa aggcctgacg ggcaagttgg acgcccga gatccgcgag 1620
attctcatta aggccaagaa gggcggcaag atcgccgtgt aataattcta ga 1672

40

15

<210> 19

<211> 1672

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 19

```
10aaagccacca tggagatgc caaaaacatt aagaagggcc cagcgccatt ctaccactc 60
   gaagacggca ccgccggcga gcagctgcac aaagccatga agcgctacgc cctgggtgcc 120
   ggcaccatcg cctttaccga cgcacatatc gaggtggaca ttacctacgc cgagtacttc 180
   gagatgagcg ttcggctggc agaagctatg aagcgctatg ggctgaacac caaccatcgc 240
   atcgtggtgt gcagcgagaa tagcttgacg ttcttcatgc ccgtggtggg tgccctgttc 300
15atcgggtgtg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360
   atgggcatca gccagcccac cgtcgtattc gtgagcaaga aagggtgca aaagatcctc 420
   aacgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
   taccagggct tccaaagcat gtacaccttc gtgacttccc atttgccacc cggcttcaac 540
   gagtacgact tcgtgcccga gagcttcgac cgggacaaaa ccatcgccct gatcatgaac 600
20agtagtggca gtaccggatt gcccaagggc gtagccctac cgcaccgcac cgcttgtgtc 660
   cgattcagtc atgcccgcga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
   ctacgcgtgg tgccatttca ccacggcttc ggcattgtca ccacgctggg ctacttgatc 780
   tgcggctttc gggctcgtgt catgtaccgc ttcgaggagg agctattctt gcgcagcttg 840
   caagactata agattcaaag cgccctgctg gtgccacac tgttcagttt cttcgccaaag 900
25agcactctca tcgacaagta cgacctaacg aacttgcacg agatcgccag cggcggggcg 960
   ccgctcagca aggaggtggg cgaggccgtg gccaaacgct tccacctacc aggcatccgc 1020
   cagggtacg gcctgacaga aacaaccagc gccattctga tcacccccga aggggacgac 1080
   aagcctggcg cagtaggcaa ggtggtgccc ttcttcgagg ctaagggtgt ggacttgga 1140
   accggaaga cactgggtgt gaaccagcgc ggcgagctgt gcgtccgtgg ccccatgatc 1200
30atgagcggct acgttaacaa ccccgaggct acaaacgctc tcatcgacaa ggacggctgg 1260
   ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
   ctgaagagcc tgatcaaata caagggttac caggtagccc cagccgaact ggagagcatc 1380
   ctgctgcaac accccaacat cttcgacgcc ggggtcgccg gcctgcccga cgacgatgcc 1440
   ggcgagctgc ccgccgagc cgtcgtgctg gaacacggta aaaccatgac cgagaaggag 1500
35atcgtggact atgtggccag ccaggttaca accgccaaga agctgcgcgg tgggtgtgtg 1560
   ttcgtggacg aggtgcctaa aggcctgacg ggcaagtgtg acgcccgcaa gatccgcgag 1620
   attctcatta aggccaagaa gggcggcaag atcgccgtgt aataattcta ga 1672
```

<210> 20

40<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

<400> 20

```
5aaagccacca tggaagatgc caaaaacatt aagaagggcc cagcgccatt ctaccctc 60
gaagacggca ccgccggcga gcagctgcac aaagccatga agcgctacgc cctgggcccc 120
ggcaccatcg cctttaccga cgcacatcgc gaggtggaca ttacctacgc cgagtacttc 180
gagatgagcg ttcggctggc agaagctatg aagcgctatg ggctgaacac caaccatcgg 240
atcgtggtgt gcagcgagaa tagcttgacg ttcttcatgc ccgtgttggg tgccctgttc 300
10atcgggtgtg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360
atgggcatca gccagccac cgctgtattc gtgagcaaga aagggtgca aaagatcctc 420
aacgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccagggct tccaaagcat gtacaccttc gtgacttccc atttgccacc cggcttcaac 540
gagtacgact tcgtgcccga gagcttcgac cgggacaaaa ccatcgccct gatcatgaac 600
15agtagtggca gtaccggatt gcccaaggcg gtagccctac cgcaccgcac cgcttgtgtc 660
cgattcagtc atgcccgcga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
ctcagcgtgg tgccatttca ccacggcttc ggcattgtca ccacgctggg ctacttgatc 780
tgcggctttc gggctcgtgt catgtaccgc ttcgaggagg agctattctt gcgcagcttg 840
caagactata agattcaaag cgccctgctg gtgccacac tgttcagttt ctctgctaag 900
20agcactctca tcgacaagta cgacctaacg aacttgacg agatcgccag cggcggggcg 960
ccgctcagca aggaggtggg cgaggccgtg gccaaacgct tccacctacc aggcattccg 1020
cagggtacg gcctgacaga aacaaccagc gccattctga tcacccccga aggggacgac 1080
aagcctggcg cagtaggcaa ggtggtgccc ttcttcgagg ctaagggtgt ggacttggac 1140
accggtaga cactgggtgt gaaccagcgc ggcgagctgt gcgtccgtgg ccccatgatc 1200
25atgagcggct acgttaacaa ccccgaggct acaaacgctc tcatcgacaa ggacggctgg 1260
ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
ctgaagagcc tgatcaaata caagggttac caggtagccc cagccgaact ggagagcatc 1380
ctgtgcaac accccaacat ctctgacgcc ggggtcgcg gcctgcccga cgacgatgcc 1440
ggcgagctgc ccgccgagt cgtcgtgctg gaacacggta aaaccatgac cgagaaggag 1500
30atcgtggact atgtggccag ccagggtaca accgccaaga agctgcgcgg tgggtgtgtg 1560
ttcgtggacg aggtgcctaa aggcctgacg ggcaagtgg acgcccga gatccgcgag 1620
attctcatta aggccaaagaa gggcggaag atcgccgtgt aataattcta ga 1672
```

<210> 21

35<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

40<223> A synthetic construct.

<400> 21

```
aaaqccacca tqqaagatgc caaaaacatt aaqaagggcc cagcgccatt ctaccctc 60
```

17

```

gaagacggca cgcggcgga gcagctgcac aaagccatga agcgctacgc cctggtgccc 120
ggcaccatcg cctttaccga cgcacatatc gaggtggaca ttacctacgc cgagtacttc 180
gagatgagcg ttcggctggc agaagctatg aagcgctatg ggctgaatac aaaccatcgg 240
atcgtggtgt gcagcgagaa tagcttgacg ttcttcatgc ccgtgttggg tgccctgttc 300
5atcgggtgtg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360
atgggcatca gccagccac cgctgtattc gtgagcaaga aagggtgca aaagatcctc 420
aacgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccagggct tccaaagcat gtacaccttc gtgacttccc atttgccacc cggttccaac 540
gagtacgact tcgtgcccga gagcttcgac cgggacaaaa ccacgcccct gatcatgaac 600
10agtagtggca gtaccggatt gcccaaggcg gtagccctac cgcaccgcac cgcttgtgtc 660
cgattcagtc atgcccgga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
ctcagcgtgg tgccatttca ccacggcttc ggcatgttca ccacgctggg ctacttgatc 780
tgcggctttc gggctcgtgt catgtaccgc ttcgaggagg agctattctt gcgcagcttg 840
caagactata agattcaaag cgccctgtg gtgccacac tgttcagttt cttcgctaag 900
15agcactctca tcgacaagta cgacctaacg aacttgacg agatcgccag cggcggggcg 960
ccgctcagca aggaggtagg tgaggccgtg gccaaacgct tccacctacc aggcattccg 1020
cagggtacg gcctgacaga aacaaccagc gccattctga tccccccga aggggacgac 1080
aagcctggcg cagtaggcaa ggtggtgccc ttcttcgagg ctaagggtgt ggacttggac 1140
accggaaga cactgggtgt gaaccagcgc ggcgagctgt gcgtccgtgg ccccatgatc 1200
20atgagcggct acgttaacaa ccccgaggct acaaacgctc tcatcgacaa ggacggctgg 1260
ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
ctgaagagcc tgatcaaata caagggtac caggtagccc cagccgaact ggagagcatc 1380
ctgtgcaac accccaacat cttcgacgcc ggggtcgccg gcctgcccga cgacgatgcc 1440
ggcgagctgc ccgccgagt cgctgtgtg gaacacggta aaaccatgac cgagaaggag 1500
25atcgtggact atgtggccag ccagggtaca accgccaaga agctgcgcgg tgggtgtgtg 1560
ttcgtggacg aggtgcctaa aggcctgacg ggcaagttgg acgcccga gatccgcgag 1620
attctcatta aggccaagaa gggcggaag atcgccgtgt aataattcta ga 1672

```

<210> 22

30<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

35<223> A synthetic construct.

<400> 22

```

aaagccacca tggaagatgc caaaaacatt aagaagggcc cagcgccatt ctaccactc 60
gaagacggga cgcggcgga gcagctgcac aaagccatga agcgctacgc cctggtgccc 120
40ggcaccatcg cctttaccga cgcacatatc gaggtggaca ttacctacgc cgagtacttc 180
gagatgagcg ttcggctggc agaagctatg aagcgctatg ggctgaatac aaaccatcgg 240
atcgtggtgt gcagcgagaa tagcttgacg ttcttcatgc ccgtgttggg tgccctgttc 300
atcgggtgtg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360

```

18

```

atgggcatca gccagccac cgctcgattc gtgagcaaga aagggtgca aaagatcctc 420
aacgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccagggct tccaaagcat gtacaccttc gtgacttccc atttgccacc cggttcaac 540
gagtacgact tcgtgcccga gagcttcgac cgggacaaaa ccatcgccct gatcatgaac 600
5agtagtggca gtaccggatt gcccaagggc gtagccctac cgcaccgcac cgcttgtgtc 660
cgattcagtc atgcccgcga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
ctcagcgtgg tgccatttca ccacggcttc ggcatgttca ccacgctggg ctacttgatc 780
tgcggtttc gggtcggtgt catgtaccgc ttcgaggagg agctattctt gcgcagcttg 840
caagactata agattcaatc tgccctgctg gtgcccacac tatttagctt cttcgctaag 900
10agcactctca tcgacaagta cgacctaac aacttgcacg agatcgccag cggcggggcg 960
ccgctcagca aggaggtagg tgaggccgtg gccaaacgct tccacctacc aggcattccg 1020
cagggtacg gcctgacaga aacaaccagc gccattctga tccccccga aggggacgac 1080
aagcctggcg cagtaggcaa ggtggtgccc ttcttcgagg ctaagggtgt ggacttgga 1140
accggtaga cactgggtgt gaaccagcgc ggcgagctgt gcgtccgtgg ccccatgatc 1200
15atgagcggct acgttaacaa ccccgaggct acaaacgctc tcatcgacaa ggacggctgg 1260
ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
ctgaagagcc tgatcaaata caagggtac caggtagccc cagccgaact ggagagcatc 1380
ctgctgcaac accccaacat cttcgacgcc ggggtcgccg gcctgcccga cgacgatgcc 1440
ggcgagctgc ccgcccagc cgctcggtgt gaacacggta aaacctgac cgagaaggag 1500
20atcggtgact atgtggccag ccagggtaca accgccaaga agctgcgcgg tgggtgtgtg 1560
ttcgtggacg aggtgcctaa aggcctgacg ggcaagtgtg acgcccga gatccgcgag 1620
attctcatta aggccaagaa gggcggaag atcgccgtgt aataattcta ga 1672

```

<210> 23

25<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

30<223> A synthetic construct.

<400> 23

```

aaagccacca tggaagatgc caaaaacatt aagaagggcc cagcgccatt ctaccactc 60
gaagacggga ccgcccgcga gcagctgcac aaagccatga agcgctacgc cctggtgccc 120
35ggcaccatcg cttttaccga cgcacatata gaggtggaca ttacctacgc cgagtacttc 180
gagatgagcg ttcggctggc agaagctatg aagcgctatg ggctgaatac aaaccatcgg 240
atcggtgtgt gcagcgagaa tagcttgacg ttcttcatgc ccgtgttggg tgccctgttc 300
atcggtgtgg ctgtggcccc agctaacgac atctacaacg agcgcgagct gctgaacagc 360
atgggcatca gccagccac cgctcgattc gtgagcaaga aagggtgca aaagatcctc 420
40aacgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccagggct tccaaagcat gtacaccttc gtgacttccc atttgccacc cggttcaac 540
gagtacgact tcgtgcccga gagcttcgac cgggacaaaa ccatcgccct gatcatgaac 600
agtagtggca gtaccggatt gcccaagggc gtagccctac cgcaccgcac cgcttgtgtc 660

```

19

```

cgattcagtc atgcccgga ccccatcttc ggcaaccaga tcatccccga caccgctatc 720
ctcagcgtgg tgccatttca ccacggttc ggcatgttca ccacgctggg ctacttgatc 780
tgcggtcttc gggtcgtgct catgtaccgc ttcgaggagg agctattctt gcgcagcttg 840
caagactata agattcaatc tgccctgctg gtgcccacac tatttagctt cttcgctaag 900
5agcactctca tcgacaagta cgacctaaag aacttgacg agatcgccag cggcggggcg 960
ccgctcagca aggaggtagg tgaggccgtg gccaaacgct tccacctacc aggcacccgc 1020
cagggctacg gcctgacaga aacaaccagc gccattctga tccccccga aggggacgac 1080
aagcctggcg cagtaggcaa ggtggtgcc ttcttcgagg ctaagggtgt ggacttgga 1140
accgtaaga cactgggtgt gaaccagcg ggcgagctgt gcgtccgtgg ccccatgatc 1200
10atgagcggct acgttaacaa ccccgaggct acaaacgctc tcatcgacaa ggacggctgg 1260
ctgcacagcg gcgacatcgc ctactgggac gaggacgagc acttcttcat cgtggaccgg 1320
ctgaagagcc tgatcaaata caagggtac caggtagccc cagccgaact ggagagcatc 1380
ctgctgcaac accccaacat cttcgacgcc ggggtcgccg gcctgcccga cgacgatgcc 1440
ggcgagctgc ccgccgcagt cgtcgtgctg gaacacggtg aaaccatgac cgagaaggag 1500
15atcgtggact atgtggccag ccagggtaca accgccaga agctgcgcgg tgggtgtgtg 1560
ttcgtggacg aggtgcctaa aggactgacc ggcaagtgg acgccgcaa gatccgcgag 1620
attctcatta aggccaaaga gggcggaag atcgccgtgt aataattcta ga 1672

```

<210> 24

20<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

25<223> A synthetic construct.

<400> 24

```

aaagccacca tggaggatgc taagaatatt aagaaggggc ctgctccttt ttatcctctg 60
gaggatggga cagctgggga gcagctgcat aaggctatga agagatatgc tctgggtgcct 120
30gggacaattg cttttacaga tgctcatatt gaggtggata ttacatatgc tgagtatttt 180
gagatgtctg tgagactggc tgaggctatg aagagatatg ggctgaatac aaatcataga 240
attgtggtgt gttctgagaa ttctctgcag ttttttatgc ctgtgctggg ggctctgttt 300
attggggtgg ctgtggctcc tgctaataatg atttataatg agagagagct gctgaattct 360
atggggattt ctacagctac agtgggtgtt gtgtctaaga aggggctgca gaagattctg 420
35aatgtgcaga agaagctgcc tattattcag aagattatta ttatggattc taagacagat 480
tatcaggggt ttcagtctat gtatacattt gtgacatctc atctgcctcc tgggtttaat 540
gagtatgatt ttgtgcctga gtcttttgat agagataaga caattgctct gattatgaat 600
tcttctgggt ctacagggct gcctaagggg gtggctctgc ctcatagaac agcttgtgtg 660
agattttctc atgctagaga tcctattttt gggaatcaga ttattcctga tacagctatt 720
40ctgtctgtgg tgccctttca tcatgggttt gggatgttta caaactggg gtatctgatt 780
tgtgggttta gagggtgct gatgtataga tttgaggagg agctgtttct gagatctctg 840
caggattata agattcagtc tgctctgctg gtgcctacac tgttttcttt ttttgctaag 900
tctacactga ttaataaata taatctatct aatctacata aatttcttc taaaaaacct 960

```

20

```

cctctgtcta aggaggtggg ggaggctgtg gctaagagat ttcattctgcc tgggattaga 1020
caggggtatg ggctgacaga gacaacatct gctattctga ttacacctga gggggatgat 1080
aagcctgggg ctgtggggaa ggtggtgcct ttttttgagg ctaaggtggg ggatctggat 1140
acagggaaga cactgggggt gaatcagaga ggggagctgt gtgtgagagg gcctatgatt 1200
5atgtctgggt atgtgaataa tcctgaggct acaaatgctc tgattgataa ggatgggtgg 1260
ctgcattctg gggatattgc ttattgggat gaggatgagc atttttttat tgtggataga 1320
ctgaagtctc tgattaagta taaggggtat cagggtggctc ctgctgagct ggagtctatt 1380
ctgctgcagc atcctaatat ttttgatgct ggggtggctg ggctgcctga tgatgatgct 1440
ggggagctgc ctgctgctgt ggtggtgctg gagcatggga agacaatgac agagaaggag 1500
10attgtggatt atgtggcttc tcaggtgaca acagctaaga agctgagagg ggggggtggg 1560
tttgtggatg aggtgcctaa ggggctgaca ggggaagctgg atgctagaaa gattagagag 1620
attctgatta aggctaagaa gggggggaag attgctgtgt aataattcta ga 1672

```

<210> 25

15<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

20<223> A synthetic construct.

<400> 25

```

aaagccacca tggaagatgc taaaaacatt aagaaggggc ctgctccttt ctaccctctg 60
gaggatggga ctgccgggga gcagctgcac aaagctatga agcggatgac tctgggtgcca 120
25ggcacaattg cgttcacgga tgctcacatt gaggtggaca ttacatacgc tgagtatttt 180
gagatgtcgg tgcggctggc tgaggctatg aagcgatatg ggctgaatac aaaccataga 240
attgtagtgt gctctgagaa ctcgttgcag ttttttatgc ctgtgctggg ggctctcttc 300
atcgggggtg ctgtggctcc tgctaacgac atttacaatg agagagagct tttgaactcg 360
atggggattt ctcagcctac agtgggtgtt gtgagtaaga aagggttca aaagattctc 420
30aatgtgcaaa agaagctgcc tattattcaa aagattatta ttatggactc taagacagac 480
taccaggggt ttcagtctat gtatacat tttgtgacatctc atctgcctcc tgggttcaac 540
gagtatgact ttgtgcccga gtctttcgac agagataaga caattgctct gattatgaat 600
tcatctgggt ctaccgggct gcctaagggt gtagctctgc cacatagaac agcttgtgtg 660
agattttctc atgctagga ccctattttt gggaatcaga ttattcctga tactgctatt 720
35ctgtcggttg tgccctttca tcatgggttt gggatgttta caacactggg ctacctgata 780
tgtgggttta gagtgggtgt catgtatagg tttgaggagg agcttttttt gcgctctctg 840
caagattata agattcagtc tgctctgctg gtgcctacac tgttttcttt ttttgctaag 900
tctaccctga tcgataagta tgatctgtcc aacctgcacg agattgcttc tgggggggct 960
cctctgtcta aggaggtagg tgaggctgtg gctaagcgtc ttcattctgcc tgggaatcaga 1020
40caggggtatg ggctaacaga aacaacatct gctattctga ttacaccaga gggggatgat 1080
aagcccgggg ctgtagggaa agtggtgccc ttttttgaag ctaaagtagt tgatcttgat 1140
accggaaga cactgggggt gaatcagcga ggggaactgt gtgtgagagg gcctatgatt 1200
atctcctcct atctcaacaa ccctcaacct acaaatgctc tgattgataa ggatgggtgg 1260

```

21

ctgcattcgg gcgatat tgc ttactgggat gaggatgagc atttcttcat cgtggacaga 1320
 ctgaagtcgt tgatcaaata taaggggtat caagtagctc ctgctgagct ggagtccatt 1380
 ctgcttcaac atcctaacat tttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
 ggggagctgc ctgctgctgt agtgggtgctg gagcacggta agacaatgac agagaaggag 1500
 5attgtggatt atgtggcttc acaagtgaca acagctaaga aactgagagg tggcgttgtg 1560
 tttgtggatg aggtgcc taa agggctgaca ggcaagctgg atgctagaaa aattcgagag 1620
 attctgatta aggctaagaa ggggtggaaag attgctgtgt aatagttcta ga 1672

<210> 26

10<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

15<223> A synthetic construct.

<400> 26

aaagccacca tggaaagtgc taaaaacatt aagaaggggc ctgctccttt ctaccctctt 60
 gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgac tcttgtgcca 120
 20ggcacaattg cgttcacgga tgctcacatt gaggtggaca tcacatacgc tgagtatttt 180
 gagatgtcgg tgcggctggc agaagctatg aagcgcctatg ggctgaatac aaaccataga 240
 attgtagtgt gcagtgaaga ctcggtgcag ttctttatgc ccgtgctggg ggctctcttc 300
 atcgggggtg ctgtggctcc tgctaacgac atctacaacg agcgagagct gttgaactcg 360
 atggggattt ctcagcctac agtgggtgtt gtgagtaaga aagggttca aaagattctc 420
 25aatgtgcaaa agaagctgcc tattattcaa aagattatta ttatggactc taagaccgac 480
 taccaggggt ttcagtctat gtatacat ttt gtgacatctc atctgcctcc tggcttcaac 540
 gagtacgact tcgtgcc cga gtctttcgac agagataaga caattgctct gatcatgaat 600
 tcatccgggt ctaccgggct gcctaagggg gtagctctgc cccatagaac agcttgtgtg 660
 agattttctc atgctaggga ccctattttt gggaatcaga ttattcctga cactgctatt 720
 30ctgtcgggtg tgccctt tca tcatgggttt gggatgttta caacactggg ctacctata 780
 tgtgggttta gagtgggtgct catgtatagg tttgaagaag agctgttctt acgctctttg 840
 caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgcataag 900
 tctacgctca tagacaagta tgacttgctc aacttgacag agattgcttc tggcggagca 960
 cctctgtcta aggaggtagg tgaggctgtg gctaagcgtc ttcactctgc tggatcaga 1020
 35caggggtatg ggctaacaga aacaacatct gctattctga ttacaccaga gggggatgat 1080
 aagccccggg ctgtagggaa agtgggtgcc ttttttgaag ccaaagtagt tgatcttgat 1140
 accggtaaga cactaggggg gaaccagcgt ggtgaactgt gtgtgagagg gcctatgatt 1200
 atgtcggggg acgttaacaa ccccgaaagc acaaatgctc tgattgataa ggatggctgg 1260
 ctgcattcgg gcgacattgc ttactgggat gaggatgagc atttcttcat cgtggacaga 1320
 40ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgctgagct ggaatccatt 1380
 ctgcttcaac atcccaacat tttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
 ggggagttgc ctgctgctgt agtgggtgct gagcacggta agacaatgac agagaaggag 1500
 atcgtggatt atgtggcttc acaagtgaca acagctaaga aactgagagg tggcgttgtg 1560

22

tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgctagaaa aattcgagag 1620
attctgatta aggctaagaa gggtagaa ag attgctgtgt aatagttcta ga 1672

<210> 27

5<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

10<223> A synthetic construct.

<400> 27

aaagccacca tggaagatgc taaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgc tcttgtgcc 120
15ggcaccaattg cgttcacgga tgctcacatt gaggtggaca tcacatacgc tgagtatttt 180
gagatgtcgg tgcggctggc agaagctatg aagcgtatg ggctgaatac aaaccataga 240
attgtagtgt gcagtgagaa ctggttcag ttctttatgc ccgtgctggg ggctctcttc 300
atcgggggtg ctgtggctcc tgctaacgac atctacaacg agcgagagct gttgaactcg 360
atggggatct ctacgcctac agtgggtgtt gtgagtaaga aagggttca aaagattctc 420
20aatgtgcaaa agaagctgcc tattattcaa aagattatta ttatggactc taagacagac 480
taccaggggt ttcagtcac gtatacat tt gtgacatctc atctgcctcc tggcttcaac 540
gagtacgact tcgtgcccga gtctttcagc agagataaga caattgctct gatcatgaat 600
tcacccgggt ctaccgggt gcctaagggt gtagctctgc cccatcgaac agcttgtgtg 660
agattctctc atgccaggga cccgatcttt gggaatcaga ttattcctga cactgctatt 720
25ctgtcgggtg tgccctttca tcatgggttt gggatgttta caacactggg atacctata 780
tgtgggttta gagtgtgtc catgtatagg ttgaagaag aactgttctt acgctctttg 840
caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgcataag 900
tctacgctca tagacaagta tgacttgtcc aacttgacg agattgcttc tggcggagca 960
cctctgtcta aggaggtagg tgaggctgtg gctaagcgt ttcatctgcc tggatcaga 1020
30caggggtacg ggctaacaga aacaacttct gctattctga ttacaccaga gggcgtgac 1080
aagccgggg ctgtaggga agtgggtgcc ttttttgaag ccaaagtagt tgatcttgat 1140
accggtaaga cactaggggt gaaccagcgt ggtgaactgt gtgtgctggg ccctatgatt 1200
atgtcgggg acgttaacaa cccgaagct acaaatgctc ttattgataa ggatggctgg 1260
ttgcattcgg gcgacattgc ctactgggat gaggatgagc atttcttcat cgtggacaga 1320
35ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgctgagct ggaatccatt 1380
ctgcttcaac atccaaacat tttcgtgctt ggggtggctg ggctgcctga tgatgatgct 1440
ggagagttgc ctgctgctgt agtagtgctt gagcacggtg agacaatgac agagaaggag 1500
atcgtggatt atgtggcttc acaagtga ca acagctaaga aactgagagg tggcgtgtg 1560
tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgccagaaa aattcgagag 1620
40attctcatta aggctaagaa gggtagaa ag attgctgtgt aatagttcta ga 1672

<210> 28

<211> 1672

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 28

```
10aaagccacca tggagatgc taaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
   gaagatggga ctgctggcga gcaacttcac aaagctatga agcggtatgc tcttgtgcca 120
   ggcacaattg cgttcacgga tgctcacatt gaggtggaca tcacatacgc tgagtatttt 180
   gagatgtcgg tgcggctggc agaagctatg aagcgcctatg ggctgaatac aaaccataga 240
   attgtagtgt gcagtgagaa ctcgttgcag ttctttatgc ccgtgctggg ggctctcttc 300
15atcgggggtg ctgtggctcc tgctaacgac atctacaacg agcgagagct gttgaactcg 360
   atggggatct ctcagcctac agtgggtgtt gtgagtaaga aagggttca aaagattctc 420
   aatgtgcaaa agaagctgcc tattatacaa aagattatta ttatggactc taagaccgac 480
   taccaggggt ttcagtccat gtacacattt gtaacctctc atctgcctcc tggcttcaac 540
   gagtacgact tcgtgcccga gtctttcgac agggacaaaa cgattgctct gatcatgaac 600
20tcatccgggt ctaccggggt gcctaagggt gtagctctgc cccatcgaa agcttggtg 660
   agattctctc atgccaggga cccgatcttt gggaaacaga ttattcctga cactgctatt 720
   ctgtcgggtg tgccctttca tcatgggttt gggatgttca caaactggg atacctcatt 780
   tgcgggttta gagtgtgtct catgtatagg tttgaagaag aactattcct acgctctttg 840
   caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgctaag 900
25tctacgctca tagacaagta tgacttgctc aacttgcacg agattgcttc tggcggagca 960
   cctctgtcta aggaggtagg tgaggctgtg gctaagcgtc ttcactctgc tgggtatcaga 1020
   caggggtacg ggctaacaga aacaacttct gctattctga ttacaccaga gggcgatgac 1080
   aaaccgggg ctgtagggaa agtgggtgcc ttttttgaag ccaaagtagt tgatcttgat 1140
   accggtaaga cactaggggt gaaccagcgt ggtgaactgt gtgtgcgggg ccctatgatt 1200
30atgtcgggg acgttaacaa cccgaagct acaaatgctc ttattgataa ggatggctgg 1260
   ttgcattcgg gcgacattgc ctactgggat gaggatgagc atttcttcat cgtggacaga 1320
   ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgctgagct ggaatccatt 1380
   ctgcttcaac atcctaacat tttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
   ggagagtgtc ctgctgctgt agtagtgctt gagcacggta agacaatgac agagaaggag 1500
35atcgtggatt atgtggcttc acaagtgaca acagctaaga aactgagagg tggcgttgtg 1560
   tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgccagaaa aattcgagag 1620
   attctcatta aggctaagaa ggggtggaaag attgctgtgt aatagttcta ga 1672
```

<210> 29

40<211> 1672

<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

<400> 29

```
5aaagccacca tggaagatgc caaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgc tcttgtgcca 120
ggcacaattg cgttcacgga tgctcacatt gaagtagaca tcacatacgc tgagtatttt 180
gagatgtcgg tgcggctggc agaagctatg aagcgtatg ggctgaatac aaaccataga 240
attgtagtgt gcagtgagaa ctcgttgagc ttctttatgc ccgtgctggg ggctctcttc 300
10atcgggggtg ctgtggctcc tgctaacgac atctacaacg agcgagagct gtggaactcg 360
atggggatct ctcagcctac agtgggtgtt gtgagtaaga aagggcttca aaagattctc 420
aatgtgcaaa agaagctgcc tattatacaa aagattatta ttatggactc taagaccgac 480
taccaggggt ttcagtcctat gtacacattt gtaacctctc atctgcctcc tggcttcaac 540
gagtacgact tcgtgcccga gtctttcgac agggacaaaa cgattgctct gatcatgaac 600
15agctccgggt ctaccgggct gcctaagggt gtagctctgc cccatcgaac agcttgtgtg 660
agattctctc atgccaggga cccgatcttt ggaaaccaga tcatccctga cactgctatt 720
ctgtcgggtg tgcccttca tcatgggttt gggatgttca caacactggg atacctcatt 780
tgcggttcta gagtgggtct catgtatagg tttgaagaag aactattcct acgctctttg 840
caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgcctaag 900
20tctacgctca tagacaagta tgacttgtcc aacttgcacg agattgcttc tggcggagca 960
cctctgtcta aggaggtagg tgaggctgtg gctaagcgct ttcactctgc tggatcaga 1020
caggggtacg ggctaacaga aacaacttct gctattctga ttacaccaga gggcgatgac 1080
aaaccggggg ctgtagggaa agtgggtgcc ttttttgaag ccaaagtagt tgatcttgat 1140
accggtgaaga cactaggggt gaaccagcgt ggtgaactgt gtgtgcgggg ccctatgatt 1200
25atgtcggggt acgttaacaa cccgaagct acaaagtctc tcatagacaa ggacgggtg 1260
cttcatagcg gcgacattgc ctactgggac gaggatgagc atttcttcat cgtggacaga 1320
ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgctgagct ggaatccatt 1380
ctgcttcaac accccaatat cttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
ggagagctgc ctgctgctgt agtagtgctt gagcacggtg agacaatgac agagaaggag 1500
30atcgtggatt atgtggcttc acaagtgaca acagctaaga aactgagagg tggcggtgtg 1560
tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgccagaaa aattcgagag 1620
attctcatta aggctaagaa gggtggaag attgctgtgt aatagttcta ga 1672
```

<210> 30

35<211> 1056

<212> DNA

<213> Artificial Sequence

<220>

40<223> A synthetic construct.

<400> 30

```
ccactcagtg gccaccatga agaagcccga gctgaccgct accagcgttg agaagttcct 60
```

25

```

gatcgagaag ttcgacagcg tgagcgacct gatgcagtta agcgagggcg aggaaagccg 120
cgccttcagc ttcgatgtcg gcggacgcgg ctatgtactg cgggtgaata gctgcgctga 180
tggcttctac aaagaccgct acgtgtaccg ccacttcgcc agcgctgcac tgcccatccc 240
cgaggtgctg gacatcggcg agttcagcga gagcctgaca tactgcatca gccgccgcgc 300
5tcaaggcgtg actctccaag acctgcccga gacagagctg cccgctgtgc tacagcctgt 360
cgccgaggct atggacgcta ttgccgccgc cgacctgagc cagaccagcg gcttcggccc 420
attcggggccc caaggcatcg gccagtacac cacctggcgc gacttcactc gcgccattgc 480
tgatcccat gtctaccact ggcagaccgt gatggacgac accgtgagcg ccagcgtagc 540
tcaagccctg gacgagctga tgctgtgggc cgaggactgc cccgaggtgc gccatctcgt 600
10ccatgccgac ttcggcagca acaacgtcct gaccgacaac ggccgcatca ccgccgtaat 660
cgactggagc gaggccatgt tcggggacag tcagtacgag gtggccaaca tcttcttctg 720
gcgcccttg ctggcctgca tggagcagca aaccgctac ttcgagcgcc gccatcccga 780
gctggccggc agcccccgtc tgcgagccta catgctgcgc atcggcctgg atcagctcta 840
ccagagcctc gtggacggca acttcgacga tgctgcctgg gctcaaggcc gctgcgatgc 900
15catcgtccgc agcggggccg gcaccgtcgg tcgcacacaa atcgtctgcc ggagcgccgc 960
cgtatggacc gacggctgcg tcgaggtgct ggccgacagc ggcaaccgcc ggcccagtac 1020
acgaccgcgc gctaaggagt agtaaccagc tcttg 1056

```

<210> 31

20<211> 1672

<212> DNA

<213> Artificial Sequence

25<220>

<223> A synthetic construct.

<400> 31

```

aaagccacca tggaagatgc caaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
30gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgc tcttggtcca 120
gggacaattg cgttcacgga tgctcacatt gaagtagaca tcacatacgc tgagtatttt 180
gagatgtcgg tcgggtggc agaagctatg aagcgtatg ggctgaatac aaaccataga 240
attgtagtgt gcagtgagaa ctggtgcag ttctttatgc ccgtgctggg ggctctcttc 300
atcgggggtg ctgtggctcc tgctaacgac atctacaacg agcgagagct gttgaactcg 360
35atggggatct ctcagcctac agtggtgttt gtgagtaaga aagggttca aaagattctc 420
aatgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccaggggt ttcagtccat gtacacatct gtaacctctc atctgcctcc tggcttcaac 540
gagtacgact tcgtgcccga gtctttcgac agggacaaaa cgattgctct gatcatgaac 600
agctccgggt ctaccggggt gcctaagggt gtagctctgc cccatcgaac agcttggtgtg 660
40agattctctc atgccaggga cccgatcttt ggaaaccaga tcatccctga cactgctatt 720
ctgtcgggtg tgccctttca tcatgggttt gggatgttca caacactggg atacctcatt 780
tgcgggttta gagtgtgtct catgtatagg tttgaagaag aactattcct acgctctttg 840
caagattata agattcaqtc tgctctgctg gtgccaacac tattctcttt ttttgctaag 900

```

```

tctacgctca tagacaagta tgacttgtcc aacttgacg agattgcttc tggcggagca 960
cctctgtcta aggaggtagg tgaggctgtg gctaagcgct ttcattctgcc tggatcaga 1020
caggggtacg ggctaacaga aacaacttct gctattctga ttacaccaga gggcgatgac 1080
aaacctgggg ctgtagggaa agtggtgccc ttttttgaag ccaaagtagt tgatcttgat 1140
5accggtaga cactaggggt gaaccagcgt ggtgaactgt gtgtgcgggg ccctatgatt 1200
atgtcgggg acgttaacaa ccccgaaagt acaaatgctc tcatagacaa ggacgggtgg 1260
cttcatagcg gcgacattgc ctactgggac gaggatgagc atttcttcat cgtggacaga 1320
ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgccgagct tgagtccatt 1380
ctgcttcaac accccaatat cttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
10ggagagctgc ctgctgctgt agtagtgctt gagcatggta agacaatgac agagaaggag 1500
atcgtggatt atgtggcttc acaagtgaac acagctaaga aactccgagg tggcgttgtg 1560
tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgccagaaa aattcgagag 1620
attctcatta aggctaagaa gggtggaag attgctgtgt aatagttcta ga 1672

```

15<210> 32

<211> 1672

<212> DNA

<213> Artificial Sequence

20<220>

<223> A synthetic construct.

<400> 32

```

aaagccacca tggaagatgc caaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
25gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgac tcttgtgcc 120
gggacaattg cgttcacgga tgctcacatt gaagtagaca tcacatacgc tgagtatttt 180
gagatgtcgg tcgggtggc agaagctatg aagcgtatg ggctgaatac aaaccataga 240
attgtagtgt gcagtgagaa ctcggtgcag ttctttatgc ccgtgctggg ggctctcttc 300
attgggggtg ctgtggctcc tgctaatac atctacaacg agcgagagct gttgaacagt 360
30atggggatct ctcagcctac agtggtgttt gtgagtaaga aagggttca aaagattctc 420
aatgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccaggggt ttcagtcctt gtacacattt gtaacctctc atctgcctcc tggcttcaat 540
gagtatgact tcgtgcccga gtctttcgac agggacaaaa cgattgctct gatcatgaac 600
agcagtggtt ctaccgggtt gcctaagggt gtagctctgc cccatcgaac agcttggtgtg 660
35agattctctc atgccaggga cccgatcttt ggaaaccaga tcatccctga cactgctatt 720
ctgtcgggtg tgccctttca tcatgggttt gggatgttca caaactggg atacctcatt 780
tcgggggtta gagggtgct catgtatagg tttgaagaag aactattcct acgctctttg 840
caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgctaag 900
tctacgctca tagacaagta tgacttgtcc aacttgacg agattgcttc tggcggagca 960
40cctctgtcta aggaggtagg tgaggctgtg gctaagcgct ttcattctgcc tggatcaga 1020
caggggtacg ggctaacaga aacaacttct gctattctga ttacaccaga gggcgatgac 1080
aaacctgggg ctgtagggaa agtggtgccc ttttttgaag ccaaagtagt tgatcttgat 1140
accggtaga cactaggggt gaaccagaga ggtgaattgt gtgtgagggg ccctatgatt 1200

```

27

```

atgtcgggggt acgttaacaa ccccgaaagt acaaagtctc tcatagacaa ggacgggtgg 1260
cttcatagtg gagatattgc ctactgggat gaagatgagc atttcttcat cgtggacaga 1320
ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgccgagct tgagtccatt 1380
ctgcttcaac accccaatat cttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
5ggagagctgc ctgctgctgt agtagtgctt gagcatggta agacaatgac agagaaggag 1500
atcgtggatt atgtggcttc acaagtgaca acagctaaga aactccgagg tggcgttggtg 1560
tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgccagaaa aattcgagag 1620
attctcatta aggctaagaa gggtggaag attgctgtgt aatagttcta ga 1672

```

10<210> 33

<211> 1672

<212> DNA

<213> Artificial Sequence

15<220>

<223> A synthetic construct.

<400> 33

```

aaagccacca tggaagatgc caaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
20gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgc tcttgtgcca 120
gggacaattg cgttcacgga tgctcacatt gaagtagaca tcacatacgc tgagtatttt 180
gagatgtcgg tcgggtggc agaagctatg aagcgtatg ggctgaatac aaaccataga 240
attgtagtgt gcagtgagaa ctcgttgca gttctttatgc ccgtgctggg ggctctcttc 300
attgggggtg ctgtggctcc tgctaatac atctacaacg agcgagagct gttgaacagt 360
25atggggatct ctcagcctac agtgggtgtt gtgagtaaga aagggttca aaagattctc 420
aatgtgcaaa agaagctacc gatcatacaa aagatcatca tcatggatag caagaccgac 480
taccaggggt ttcagtcctat gtacacattt gtaacctctc atctgcctcc tggcttcaat 540
gagtatgact tcgtgcccga gtctttcgac agggacaaaa cgattgctct gatcatgaac 600
agcagtgggt ctaccgggt gcctaagggt gtagctctgc cccatcgaac agcttgtgtg 660
30agattctctc atgccaggga cccgatcttt ggaaaccaga tcatccctga cactgctatt 720
ctgtcgggtg tgccctttca tcatgggttt gggatgttca caacactggg atacctcatt 780
tcgggggtta gagtggtgct :catgtatagg tttgaagaag aactattcct acgctctttg 840
caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgctaag 900
tctacgctca tagacaagta tgacttgtcc aacttgacg agattgcttc tggcggagca 960
35cctctgtcta aggaggtagg tgaggctgtg gctaagcgtt tcatctgcc tggatatcaga 1020
caggggtacg ggctaacaga aacaacttct gctattctga ttacaccaga gggcgatgac 1080
aaacctgggg ctgtaggga agtgggtgcc ttttttgaag ccaaagtagt tgatcttgat 1140
accggaaga cactaggggt gaaccagaga ggtgaattgt gtgtgagggg ccctatgatt 1200
atgtcgggggt acgttaacaa ccccgaaagt acaaagtctc tcatagacaa ggacgggtgg 1260
40cttcatagtg gagatattgc ctactgggat gaagatgagc atttcttcat cgtggacaga 1320
ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgccgagct tgagtccatt 1380
ctgcttcaac accccaatat cttcgatgct ggggtggctg ggctgcctga tgatgatgct 1440
ggagagctgc ctgctgctgt agtagtgctt gagcatggta agacaatgac agagaaggag 1500

```

28

atcgtggatt atgtggcttc acaagtgaca acagctaaga aactccgagg tggcgttgtg 1560
tttgtggatg aggtgcctaa aggactcact ggcaagctgg atgccagaaa aattcgagag 1620
attctcatta aggctaagaa gggtaggaaag attgctgtgt aatagttcta ga 1672

5<210> 34

<211> 10

<212> DNA

<213> Artificial Sequence

10<220>

<223> A synthetic construct.

<400> 34

gccaccatga

10

15

<210> 35

<211> 11

<212> DNA

<213> Artificial Sequence

20

<220>

<223> A synthetic construct.

<220>

25<221> misc_feature

<222> 4, 5, 6, 7, 8

<223> n = A,T,C or G

<400> 35

30ccannnnntg g

11

<210> 36

<211> 25

<212> DNA

35<213> Artificial Sequence

<220>

<223> A synthetic construct.

40<220>

<221> misc_feature

<222> 1, 2, 3, 4, 5, 9, 10, 11, 12, 13

<223> n = A,T,C or G

<400> 36

nnnnnccann nnntggccac catgg

25

<210> 37

5<211> 20

<212> DNA

<213> Artificial Sequence

<220>

10<223> A synthetic construct.

<220>

<221> misc_feature

<222> 10, 11, 12, 13, 14, 18, 19, 20

15<223> n = A,T,C or G

<400> 37

taataaccan nnnntggnnn

20

20<210> 38

<211> 825

<212> DNA

<213> Artificial Sequence

25<220>

<223> A synthetic construct.

<400> 38

ccactcagtg gccaccatga tgcagcagga cggcctccat gctggcagtc ccgcagcctg 60
30ggtcgagcgc ttgttcgggt acgactgggc ccagcagacc atcggatgta gcgatgccgc 120
agtgttccgc ctgagcgtc aaggccggcc cgtgctgttc gtgaagaccg acctgagcgg 180
cgccctgaac gagcttcaag acgaggctgc ccgcctgagc tggctggcca ccaccggtgt 240
accctgcgcc gctgtgttg atgttgtgac cgaagccggc cgcgactggc tgctgtggg 300
cgaggtgcct ggccaggacc tgctgagcag ccacctggcc cccgctgaga aggtgagcat 360
35catggccgac gccatgcggc gcctgcacac cctggacccc gctacatgcc ctttcgacca 420
ccaggctaag caccgcatcg agcgggctcg gaccgcgatg gaggccggcc tgggtggacca 480
ggacgacctg gacgaggagc accagggcct ggccccgct gaactgttcg cccgcctgaa 540
agccgcgatg ccggacggtg aggacctggt tgtgacacac ggcgacgct gcctccctaa 600
catcatggtc gagaacgggc gcttctccgg cttcatcgac tgcggccgcc tgggcgttgc 660
40cgaccgctac caggacatcg ccctggccac ccgcgacatc gccgaggagc tgggcggcga 720
gtgggcccgc cgcttctctg tcttgtagcg catcgagct cccgacagcc agcgcacgc 780
cttctaccgc ctgctggacg agttcttcta gtaaccaggc tctgg 825

30

<210> 39

<211> 825

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 39

```
10ccactccgtg gccacatga tcgaacaaga cggcctccat gctggcagtc ccgcagcttg 60
   ggtcgaacgc ttgttcgggt acgactgggc ccagcagacc atcggatgta gcgatgcggc 120
   cgtgttcctg ctaagcgctc aaggccggcc cgtgctgttc gtgaagaccg acctgagcgg 180
   cgccctgaac gagcttcaag acgaggctgc ccgcctgagc tggctggcca ccaccgggtg 240
   accctgcgcc gctgtgttgg atgttgtgac cgaagccggc cgggactggc tgctgctggg 300
15cgagggtccct ggccaggatc tgctgagcag ccaccttgcc cccgctgaga aggtttccat 360
   catggccgat gcaatgcggc gcctgcacac cctggacccc gctacatgcc ctttcgacca 420
   ccaggctaag catcggatcg agcgtgctcg gaccgcgatg gaggccggcc tgggtggacca 480
   ggacgacctg gacgaggagc atcagggcct gggccccgct gaactgttcg cccgcctgaa 540
   agcccgcatg ccggacgggtg aggacctggt tgtgacacat ggagatgcct gcctccctaa 600
20catcatggtc gagaatggcc gcttctccgg cttcatcgac tgcggtcgcc taggagttgc 660
   cgaccgctac caggacatcg ccctggccac ccgcgacatc gctgaggagc ttggcggcga 720
   gtgggcccgc cgcttcttag tcttgtacgg catcgcagct cccgacagcc agcgcacgcg 780
   cttctaccgc ctgctcgacg agttctttta atgaccaggc tctgg 825
```

25<210> 40

<400> 40

000

30<210> 41

<211> 861

<212> DNA

<213> Escherichia coli

35<400> 41

```
atgagtattc aacatttccg tgtcgccctt attccctttt ttgcggcatt ttgccttcc 60
gttttttgct acccagaaac gctggtgaaa gtaaaagatg ctgaagatca gttgggtgca 120
cgagtgggtt acatcgaact ggatctcaac agcggtaaga tccttgagag ttttcgcccc 180
gaagaacgtt ttccaatgat gagcactttt aaagttctgc tatgtggcgc ggtattatcc 240
40cgtattgacg ccgggcaaga gcaactcggc cgccgcatac actattctca gaatgacttg 300
gttgagtact caccagtcac agaaaagcat cttacggatg gcatgacagt aagagaatta 360
tgcagtgctg ccataaccat gagtgataac actgcggcca acttacttct gacaacgata 420
ggaggaccga aggagctaac cgcttttttg cacaacatgg gggatcatgt aactcgctt 480
```


31

gatcgttggg aaccggagct gaatgaagcc ataccaaacg acgagcgtga caccacgatg 540
cctgtagcaa tggcaacaac gttgcgcaaa ctattaactg gcgaactact tactctagct 600
tcccggcaac aattaataga ctggatggag gcggataaag ttgcaggacc acttctgcgc 660
tcggcccttc cggctggctg gtttattgct gataaatctg gagccggtga gcgtgggtct 720
5cgcggtatca ttgcagcact ggggccagat ggtaagccct cccgtatcgt agttatctac 780
acgacgggga gtcaggcaac tatggatgaa cgaaatagac agatcgctga gatagggtgcc 840
tcactgatta agcattggta a 861

<210> 42

10<211> 1056

<212> DNA

<213> Artificial Sequence

<220>

15<223> A synthetic construct.

<400> 42

ccactcctg gccacatga agaagcccga gctgaccgct accagcgttg aaaaatttct 60
catcgagaag ttgcagactg tgagcgacct gatgcagttg tcggagggcg aagagagccg 120
20agccttcagc ttcatgtcg gcggacgcgg ctatgtactg cgggtgaata gctgcgctga 180
tggttctac aaagaccgct acgtgtaccg ccacttcgcc agcgtgcac taccatccc 240
cgaagtgtg gacatcggcg agttcagcga gagcctgaca tactgcatca gtagacgcgc 300
ccaaggcgtt actctccaag acctccccga aacagagctg cctgctgtgt tacagcctgt 360
cgccgaagct atggatgcta ttgccgccgc cgacctcagt caaaccagcg gcttcggccc 420
25attcgggccc caaggcatcg gccagtacac aacctggcgg gatttcattt gcgccattgc 480
tgatcccat gtctaccact ggcagaccgt gatggacgac accgtgtccg ccagcgtagc 540
tcaagccctg gacgaactga tgctgtgggc cgaagactgt cccgaggtgc gccacctcgt 600
ccatgccgac ttccgacga acaacgtcct gaccgacaac ggccgcatca ccgccgtaat 660
cgactggagc gaggtatgt tcggggacag tcagtacgag gtggccaaca tcttcttctg 720
30gcggccctgg ctggcttgca tggagcagca gactcgctac ttcgagcgcc ggcacccga 780
gctggccggc agccctcgtc tgcgagccta catgctgcgc atcggcctgg atcagctcta 840
ccagagcctc gtggacggca acttcgacga tgctgcctgg gctcaaggcc gctgcgatgc 900
catcgtccgc agcggggccg gcaccgtcgg tcgcacacaa atcgctcgcc ggagcgccgc 960
cgtatggacc gacggctgcg tcgaggtgct ggccgacagc ggcaaccgcc ggcccagtac 1020
35acgaccgcgc gctaaggagt agtaaccagc tcttgg 1056

<210> 43

<211> 1653

<212> DNA

40<213> Artificial Sequence

<220>

<223> A synthetic construct.

<400> 43
 atggaagacg ccaaaaacat aaagaaaggc cggcgccat tctatccgct ggaagatgga 60
 accgctggag agcaactgca taaggctatg aagagatacg ccctgggtcc tggaacaatt 120
 gcttttacag atgcacatat cgagggtggac atcacttacg ctgagtactt cgaaatgtcc 180
 5gttcgggttg cagaagctat gaaacgatat gggctgaata caaatcacag aatcgctgta 240
 tgcagtgaag actctcttca attctttatg ccggtgttg ggcggttatt tatcgagatt 300
 gcagttgctc ccgcgaacga cttttataat gaacgtgaat tgctcaacag tatgggcatt 360
 tcgcagccta ccgtgggtgt cgtttccaaa aaggggttgc aaaaaatttt gaacgtgcaa 420
 aaaaagctcc caatcatcca aaaaattatt atcatggatt ctaaaacgga ttaccaggga 480
 10tttcagtcga tgtacacgtt cgtcacatct catctacctc ccggttttaa tgaatacgat 540
 tttgtgccag agtccttcga tagggacaag acaattgcac tgatcatgaa ctctcttga 600
 tctactggtc tgccataaagg tgcgctctg cctcatagaa ctgcctgcgt gagattctcg 660
 catgccagag atcctatctt tggcaatcaa atcattccgg atactgcgat ttaagtgtt 720
 gttccattcc atcacgggtt tggaatgttt actacactcg gatatttgat atgtggattt 780
 15cgagtcgtct taatgtatag atttgaagaa gagctgtttc tgaggagcct tcaggattac 840
 aagattcaaa gtgcgtgct ggtgccaacc ctattctct tcttcgcaa aagcactctg 900
 attgacaaat acgatttatc taatttacac gaaattgctt ctgggtggcg tccccctct 960
 aaggaagtcg gggaagcggc tgccaagagg ttccatctgc caggatcag gcaaggatat 1020
 gggctcactg agactacatc agctattctg attacaccgc agggggatga taaaccgggc 1080
 20gcggtcggta aagttgttcc attttttgaa gcgaagggtt tggatctgga taccgggaaa 1140
 acgctgggag ttaatcaaag aggcgaactg tgtgtgagag gtcctatgat tatgtccgg 1200
 tatgtaaaca atccggaagc gaccaacgcc ttgattgaca aggatggatg gctacattct 1260
 ggagacatag cttactggga cgaagacgaa cacttcttca tcgttgaccg cctgaagtct 1320
 ctgattaagt acaaaggcta tcaggtggct cccgctgaat tggaatccat cttgctccaa 1380
 25caccccaaca tcttcgacgc aggtgtcgca ggtcttcccg acgatgacgc cgggtgaactt 1440
 cccgccgccg ttgttgtttt ggagcacgga aagacgatga cggaaaaaga gatcgtggat 1500
 tacgtcgcca gtcaagtaac aaccgcgaaa aagttgcgcg gaggagtgtt gtttgtggac 1560
 gaagtaccga aaggtcttac cggaaaactc gacgcaagaa aaatcagaga gatcctcata 1620
 aaggccaaga agggcggaag gatcgccgtg taa 1653

30

<210> 44

<211> 1369

<212> DNA

35<213> Artificial Sequence

<220>

<223> A synthetic construct.

40<400> 44

ggatccgttt gcgtattggg cgctcttccg ctgatctgcg cagcaccatg gcctgaaata 60
 acctctgaaa gaggaacttg gttagctacc ttctgaggcg gaaagaacca gctgtggaat 120
 qtatqtcaqt taqqqtatqg aaagtcccca qqctcccaq caqqcaqaq tatqcaaaqc 180

33

```

atgcatctca attagtcagc aaccaggtgt ggaaagtccc caggctcccc agcaggcaga 240
agtatgcaaa gcatgcatct caattagtca gcaaccatag tcccgccctt aactccgccc 300
atcccccccc taactccgcc cagttccgcc cattctccgc cccatggctg actaattttt 360
tttatttatg cagaggccga ggccgcctct gcctctgagc tattccagaa gtagtgagga 420
5ggcttttttg gaggcctagg cttttgcaaa aagctcgatt cttctgacac tagcgccacc 480
atgatcgaac aagacggcct ccatgctggc agtcccgag cttgggtcga acgcttggtc 540
gggtacgact gggcccagca gaccatcgga tgtagcgatg cggccgtggt ccgtctaagc 600
gctcaaggcc ggcccgctgt gtctgtgaag accgacctga gcggcgccct gaacgagctt 660
caagacgagg ctgccgcctt gagctggctg gccaccaccg gcgtaccctg cgccgctgtg 720
10ttggatgttg tgaccgaagc cggccgggac tggctgctgc tgggcgaggt ccctggccag 780
gatctgctga gcagccacct tgccccgctt gagaagggtt ctatcatggc cgatgcaatg 840
cggcgcttgc acaccctgga ccccgctacc tgcccttctg accaccaggc taagcatcgg 900
atcgagcgtg ctccggaccg catggaggcc ggccctgggtg accaggacga cctggacgag 960
gagcatcagg gcctggcccc cgctgaactg ttcccccagc tgaaagcccg catgccggac 1020
15ggtgaggacc tggttgtcac acacggagat gcctgcctcc ctaacatcat ggtcgagaat 1080
ggccgcttct ccggcttcat cgactgcggt cgcctaggag ttgccgaccg ctaccaggac 1140
atcgccctgg ccacccgcga catcgctgag gagcttggcg gcgagtgggc cgaccgcttc 1200
ttagtcttgt acggcatcgc agctcccgac agccagcgca tcgccttcta ccgcttgctc 1260
gacgagttct tttaatgac tagaaccggt catggccgca ataaaatatc tttattttca 1320
20ttacatctgt gtgttggttt tttgtgtgtt cgaactagat gctgtcgac 1369

```

<210> 45

<211> 1214

<212> DNA

25<213> Artificial Sequence

<220>

<223> A synthetic construct.

30<400> 45

```

gcggccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
ctcagcgatc tgctatttcc gtctgtccat agtggcctga ctcccgcctg tgtagatcac 120
tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
ttcaccggcc cccgatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
35tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
aagaagtctc ccagtgaagta gtttccgaag agttgtggcc attgtacttg gcatcggtgt 360
atcacgctcg tcgttcggta tggcttcggt caactctggt tcccagcggc caagccgggt 420
cacatgatca cccatattat gaagaaatgc agtcagctcc ttagggcctc cgatcggtgt 480
cagaagtaag ttggccgcgg tggtgtcgct catggtaatg gcagcactac acaattctct 540
40taccgtcatg ccatccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtcggt 600
ttgtgagtag tgtatacggc gaccaagctg ctcttgcccc gcgtctatac gggacaacac 660
cgcgccacat agcagtactt tgaaagtgtc catcatcggg aatcgttctt cggggcgga 720
agactcaaag atcttgccgc tattgagatc cagttcgata tagccactc ttgcaccag 780

```

34

```

ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
aaatgccgca aagaagggaa tgagtgcgac acgaaaatgt tggatgctca tactcttcct 900
ttttcaatat gtttgcagca tttgtcaggg ttactagtac gtctctcttg agagaccgcg 960
atcgccacca tgtctaggta ggtagtaaac gaaagggcctt aaaggcctaa gtggccctcg 1020
5agtccagcct tgagttgggt gagtccaagt cacgtttgga gatctggtag cttacgcgta 1080
tgaggggtga gtccaagtca cgtttggaga tctggtacct tacgcgtatg agctctacgt 1140
agctagcggc ctcgggcgcc gaattcttgc gttcgaagct tggcaatccg gtactgttgg 1200
taaagccacc atgg                                     1214

```

10<210> 46

<211> 1522

<212> DNA

<213> Artificial Sequence

15<220>

<223> A synthetic construct.

<400> 46

```

gcggccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
20ctcagcgatc tgcctatttc gttcgtccat agtggcctga ctcccgtcg tgtagatcac 120
tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
ttcaccggcc cccgatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
aagaagtctg ccagttagta gtttcggaag agttgtggcc attgctactg gcatcgttgt 360
25atcacgctcg tcgttcggta tggcttcgtt caactctggt tcccagcggg caagccgggt 420
cacatgatca cccatattat gaagaaatgc agtcagctcc ttagggcctc cgatcgttgt 480
cagaagtaag ttggccgagg tggtgtcgtc catggtaatg gcagcactac acaattctct 540
taccgtcatg ccatccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtcgtt 600
ttgtgagtag tgtatacggc gaccaagctg ctcttgcccg gcgtctatac gggacaacac 660
30cgcgccacat agcagtactt tgaaagtgct catcatcggg aatcgttctt cggggcgga 720
agactcaagg atcttgccgc tattgagatc cagttcgata tagccactc ttgcaccag 780
ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
aaatgccgca aagaagggaa tgagtgcgac acgaaaatgt tggatgctca tactcttcct 900
ttttcaatat gtttgcagca tttgtcaggg ttactagtac gtctctcaag agatttgtgc 960
35atacacagtg actcatactt tcaccaatac tttgcatttt ggataaatac tagacaactt 1020
tagaagtga tttttatga ggttgtctta aaattaaaaa ttacaagta ataaatcaca 1080
ttgtaagtga ttttgtgtga taccagagg tttaaggcaa cctattactc ttatgctcct 1140
gaagtccaca attcacagtc ctgaactata atcttatctt tgtgattgct gagcaaattt 1200
gcagtataat ttcagtgtt ttaattttg tctgtcttac tatttctctt ttttatttgg 1260
40gtttgatatg cgtgcacaga atggggcttc tattaataa ttcttgagag accgcgatcg 1320
ccaccatgtc taggtaggta gtaaacgaaa gggcttaag gcctaagtgg ccctcgagtc 1380
cagccttgag ttggttgagt ccaagtcacg tttggagatc tggtaacctt cgcgtatgag 1440
ctctacgtag ctagcgccct cggcgccga attcttgcgt tcgaagcttg gcaatccggg 1500

```

actgttggtgta aagccaccat gg

1522

<210> 47

<211> 1134

5<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

10

<400> 47

gcggccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
ctcagcgatc tgcctatttc gttcgtccat agtggcctga ctcccgcgcg tgtagatcac 120
tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
15ttcaccggcc cccgatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
aagaagtctg ccagtgagta gtttcgaag agttgtggcc attgctactg gcatcgtggt 360
atcacgctcg tcgttcggta tggcttcgtt caactctggt tcccagcggc caagccgggt 420
cacatgatca cccatattat gaagaaatgc agtcagctcc ttagggcctc cgatcgttgt 480
20cagaagtaag ttggccgcgg tgttgctgct catggtaatg gcagcactac acaattctct 540
tacggtcatg ccatccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtcgtt 600
ttgtgagtag tgtatacggc gaccaagctg ctcttgcccc gcgtctatac gggacaacac 660
cgcgccacat agcagtactt tgaaagtgct catcatcggg aatcgttctt cggggcggaa 720
agactcaagg atcttgccgc tattgagatc cagttcgata tagccactc ttgcaccag 780
25ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
aaatgccgca aagaaggga tgagtgcgac acgaaaatgt tggatgctca tactcgtcct 900
ttttcaatat tattgaagca tttatcaggg ttactagtac gtctctcaag agatttgtgc 960
atacacagtg actcatactt tcaccaatac tttgcatttt ggataaatac tagacaactt 1020
tagaagtga ttatttatga gggtgtctta aaattaaaaa ttacaaagta ataatcaca 1080
30ttgtaatgta ttttgtgtga taccagagg tttaaggcaa cctattactc ttat 1134

<210> 48

<211> 319

<212> DNA

35<213> Artificial Sequence

<220>

<223> A synthetic construct.

40<400> 48

actagtacgt ctctcaagga taagtaagta atattaaggt acgggaggta cttggagcgg 60
ccgcaataaa atatctttat tttcattaca tctgtgtgtt gggtttttgt gtgaatcgat 120
aqtactaaca tacqctctcc atcaaaacaa aacqaaacaa aacaaactaq caaaataaqqc 180

36

tgtccccagt gcaagtgcag gtgccagaac atttctctgg cctaagtggc cggtaccgag 240
ctcgctagcc tcgaggatat cagatctggc ctcggcggcc aagcttggca atccgggtact 300
gttggttaaag ccaccatgg 319

5<210> 49

<211> 320

<212> DNA

<213> Artificial Sequence

10<220>

<223> A synthetic construct.

<400> 49

actagtacgt ctctcaagga taagtaagta atattaaggt acgggaggta ttggacaggc 60
15cgcaataaaa tatctttatt ttcattacat ctgtgtgttg gttttttgtg tgaatcgata 120
gtactaacat acgctctcca tcaaaacaaa acgaaacaaa acaaactagc aaaataggct 180
gtccccagtg caagtgcagg tgccagaaca tttctctggc ctaactggcc ggtacctgag 240
ctcgctagcc tcgaggatat caagatctgg cctcggcggc caagcttggc aatccgggtac 300
gttggtgtaa gccaccatgg 320

20

<210> 50

<211> 5

<212> DNA

<213> Artificial Sequence

25

<220>

<223> A synthetic construct.

<400> 50

30tataa

5

<210> 51

<211> 6

<212> DNA

35<213> Artificial Sequence

<220>

<223> A synthetic construct.

40<400> 51

stratg

6

<210> 52
<211> 9
<212> DNA
5<213> Artificial Sequence

<220>
<223> A synthetic construct.

10<220>
 <221> misc_feature
 <222> 4, 6, 7
 <223> n = A,T,C or G

15<400> 52
 mttncnnma

9

 <210> 53
 <211> 5
20<212> DNA
 <213> Artificial Sequence

 <220>
 <223> A synthetic construct.

25
 <400> 53
 tratg

5

 <210> 54
30<211> 38
 <212> DNA
 <213> Artificial Sequence

 <220>
35<223> A synthetic construct.

 <400> 54
 gtactgagac gacgccagcc caagcttagg cctgagtg

38

40<210> 55
 <211> 38
 <212> DNA
 <213> Artificial Sequence

38

<220>

<223> A synthetic construct.

<400> 55

5ggcatgagcg tgaactgact gaactagcgg ccgccgag

38

<210> 56

<211> 24

<212> DNA

10<213> Artificial Sequence

<220>

<223> A synthetic construct.

15<400> 56

ggatcccatg gtgaagcgtg agaa

24

<210> 57

<211> 21

20<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

25

<400> 57

ggatcccatg gtgaaacgcg a

21

<210> 58

30<211> 31

<212> DNA

<213> Artificial Sequence

<220>

35<223> A synthetic construct.

<400> 58

ctagcttttt tttctagata atcatgaaga c

31

40<210> 59

<211> 32

<212> DNA

<213> Artificial Sequence

39

<220>

<223> A synthetic construct.

<400> 59

5gcgtagccat ggtaaagcgt gagaaaaatg tc

32

<210> 60

<211> 33

<212> DNA

10<213> Artificial Sequence

<220>

<223> A synthetic construct.

15<400> 60

ccgactctag attactaacc gccggccttc acc

33

<210> 61

<211> 54

20<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

25

<400> 61

caaaaagcctt ggcattccgg tactgttggt aaagccacca tggatgaagcg agag

54

<210> 62

30<211> 26

<212> DNA

<213> Artificial Sequence

<220>

35<223> A synthetic construct.

<400> 62

caattgttgt tgtaacttg tttatt

26

40<210> 63

<400> 63

000

40

<210> 64

<400> 64

000

5

<210> 65

<211> 10

<212> DNA

<213> Artificial Sequence

10

<220>

<223> A synthetic construct.

<400> 65

15caccatggct

10

<210> 66

<211> 40

<212> DNA

20<213> Artificial Sequence

<220>

<223> A synthetic construct.

25<400> 66

aaccatggct tccaaggtgt acgaccccgga gcaacgcaaa

40

<210> 67

<211> 40

30<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

35

<400> 67

gctctagaat tactgctcgt tcttcagcac gcgctccacg

40

<210> 68

40<211> 31

<212> DNA

<213> Artificial Sequence

41

<220>

<223> A synthetic construct.

<400> 68

5cgctagccat ggcttcgaaa gtttatgatc c

31

<210> 69

<211> 25

<212> DNA

10<213> Artificial Sequence

<220>

<223> A synthetic construct.

15<400> 69

ggccagtaac tctagaatta ttggt

25

<210> 70

<211> 1092

20<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

25

<400> 70

aagcttgcta gcgccaccat gaagaagccc gagctcaccg ctaccagcgt tgaaaaattt 60
ctcatcgaga agttcgacag tgtgagcgac ctgatgcagt tgcggaggag cgaagagagc 120
cgagccttca gcttcgatgt cggcggagcg ggctatgtac tgcgggtgaa tagctgcgct 180
30gatggccttct acaaagaccg ctacgtgtac cgccacttcg ccagcgctgc actaccatc 240
cccgaagtgt tggacatcgg cgagttcagc gagagcctga catactgcat cagtagacgc 300
gccaagggcg ttactctcca agacctcccc gaaacagagc tgectgctgt gttacagcct 360
gtcgccgaag ctatggatgc tattgccgcc gccgacctca gtcaaaccag cggcttcggc 420
ccattcgggc cccaaggcat cggccagtac acaacctggc gggatttcat ttgcgccatt 480
35gctgatcccc atgtctacca ctggcagacc gtgatggacg acaccgtgtc cgccagcgta 540
gctcaagccc tggacgaact gatgctgtgg gccgaagact gtcccagggt gcgccacctc 600
gtccatgccg acttcggcag caacaacgtc ctgaccgaca acggccgcat caccgccgta 660
atcgactggg ccgaagctat gttcggggac agtcagtacg aggtggccaa catcttcttc 720
tgggcgccct ggctggcttg catggagcag cagactcgct acttcgagcg ccggcatccc 780
40gagctggccg gcagccctcg tctgcgagcc tacatgctgc gcatcgccct ggatcagctc 840
taccagagcc tcgtggacgg caacttcgac gatgtcgctt gggctcaagg ccgctgcgat 900
gccatcgctc gcagcggggc cggcacgctc ggtcgcacac aaatcgctcg ccggagcgcc 960
qccgtatgga ccgacggctg cgtcgagggt ctggccgaca gcggcaaccg ccggcccagt 1020

42

acacga~~c~~cg~~c~~ gcgctaagga ggg~~t~~ggcgga gggagcggtg g~~c~~ggagg~~t~~tc ctacgtatag 1080
tctagactcg ag 1092

<210> 71

5<211> 1093

<212> DNA

<213> Artificial Sequence

<220>

10<223> A synthetic construct.

<400> 71

aagcttg~~c~~ta ggc~~c~~accat gaagaagccc gagctcaccg ctaccagcgt tgaaaaattt 60
ctcatc~~g~~aga agttcgacag tgtgagcgac ctgatgcagt tgtcggagg~~g~~ cgaagagagc 120
15cgagcct~~t~~tca gcttcgatgt cggcgga~~c~~gc ggctatgtac tg~~c~~gggtgaa tagctgcgct 180
gatggc~~t~~tct acaaagaccg ctacgtgtac cgccacttcg ccagcgc~~t~~gc actaccatc 240
cccgaag~~t~~gt tggacatcgg cgagttcagc gagagcctga catactgcat cagtagacgc 300
gccc~~a~~aggcg ttactctcca agacctcccc gaaacagagc tgcctgctgt gttacagcct 360
gtcgc~~c~~gaag ctatggatgc tattgcccgc gccgacctca gtcaaaccag cggcttcggc 420
20ccattc~~g~~ggc cccaaggcat cggccagtac acaacctggc gggatttcat ttgcccatt 480
gctgat~~c~~ccc atgtctacca ctggcagacc gtgatggacg acaccgtgtc cgccagcgta 540
gctcaag~~c~~ccc tggacgaact gatgctgtgg gccgaagact gtcccagggt gcgccacctc 600
gtccat~~c~~ccg acttcggcag caacaacgtc ctgaccgaca acggccgcat caccgcgta 660
atcgac~~t~~gg~~t~~ ccgaagctat gttcggggac agtcagtacg aggtggccaa catcttcttc 720
25tggcgg~~c~~cct ggctggcttg catggagcag cagactcgt acttcgagcg ccggcatccc 780
gagctg~~g~~ccg gcagccctcg tctgcgagcc tacatgctgc gcacggcct ggatcagctc 840
taccagag~~c~~c tcgtggacgg caacttcgac gatgctgcct gggctcaagg ccgctgcgat 900
gccatc~~g~~tcc gcagcggggc cggcaccgtc ggtcgcacac aaatcgctcg ccggagcgca 960
gccgtat~~g~~ga ccgacggctg cgtcgagg~~t~~g ctggccgaca gcggcaaccg ccggcccagt 1020
30acacga~~c~~cg~~c~~ gcgctaagga aggcggtgga ggtagtgtg g~~c~~ggaggtag ctacgtataa 1080
ctctagactc gag 1093

<210> 72

<211> 813

35<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

40

<400> 72

gctagc~~g~~cca ccatgatcga acaagacggc ctccatgctg gcag~~t~~ccccg agcttggg~~t~~c 60
gaacgc~~t~~tgt tcgggtacga ctgggcccag cagaccatcg gatgtagcga tg~~c~~ggccgtg 120

43

```
ttccgtctaa gcgctcaagg ccggcccgtg ctgttcgtga agaccgacct gagcggcgcc 180
ctgaacgagc ttcaagacga ggctgcccgc ctgagctggc tggccaccac cgggtgtacc 240
tgcgccgctg tggttgatgt tgtgaccgaa gccggccggg actggctgct gctgggcgag 300
gtccctggcc aggatctgct gagcagccac cttgcccccg ctgagaaggt ttccatcatg 360
5gccgatgcaa tgcggcgccct gcacaccctg gaccccgcta catgcccctt cgaccaccag 420
gctaagcatc ggatcgagcg tgctcggacc cgcattggagg ccggcctggg ggaccaggac 480
gacctggacg aggagcatca gggcctggcc cccgctgaac tggtcgcccg cctgaaagcc 540
cgcatgccgg acggtgagga cctggttgtg acacatgggt atgcctgcct ccctaacatc 600
atggtcgaga atggccgctt ctccggcttc atcgactgcg gtcgcctagg agttgccgac 660
10cgctaccagg acatcgccct ggccaccgcg gacatcgctg aggagcttgg cggcgagtgg 720
gccgaccgct tcttagtctt gtacggcatc gcagctcccg acagccagcg catcgccctt 780
taccgctgc tcgacgagtt cttttaatct aga 813
```

<210> 73

15<211> 816

<212> DNA

<213> Artificial Sequence

<220>

20<223> A synthetic construct.

<400> 73

```
gctagcgcca ccatgatcga acaagacggc ctccatgctg gcagtcgccg agcttgggtc 60
gaacgcttgt tcgggtacga ctgggcccag cagaccatcg gatgtagcga tgcggccgtg 120
25ttccgtctaa gcgctcaagg ccggcccgtg ctgttcgtga agaccgacct gagcggcgcc 180
ctgaacgagc ttcaagacga ggctgcccgc ctgagctggc tggccaccac cggcgtacc 240
tgcgccgctg tggttgatgt tgtgaccgaa gccggccggg actggctgct gctgggcgag 300
gtccctggcc aggatctgct gagcagccac cttgcccccg ctgagaaggt ttctatcatg 360
gccgatgcaa tgcggcgccct gcacaccctg gaccccgcta cctgcccctt cgaccaccag 420
30gctaagcatc ggatcgagcg tgctcggacc cgcattggagg ccggcctggg ggaccaggac 480
gacctggacg aggagcatca gggcctggcc cccgctgaac tggtcgcccg actgaaagcc 540
cgcatgccgg acggtgagga cctggttgtc acacacggag atgcctgcct ccctaacatc 600
atggtcgaga atggccgctt ctccggcttc atcgactgcg gtcgcctagg agttgccgac 660
cgctaccagg acatcgccct ggccaccgcg gacatcgctg aggagcttgg cggcgagtgg 720
35gccgaccgct tcttagtctt gtacggcatc gcagctcccg acagccagcg catcgccctt 780
taccgcttgc tcgacgagtt cttttaatga tctaga 816
```

<210> 74

<211> 1252

40<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

<400> 74

```
5gcgggccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
ctcagcgatc tgcctatttc gttcgtccat agtggcctga ctccccgtcg tgtagatcac 120
tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
ttcaccggcc cccgatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
10aagaagtctg ccagtgagta gtttccgaag agttgtggcc attgctactg gcatcgtggg 360
atcacgctcg tcgttcggta tggcttcgtt caactctggt tcccagcggg caagccgggt 420
cacatgatca cccatattat gaagaaatgc agtcagctcc ttagggcctc cgatcgttgt 480
cagaagtaag ttggcccgcg tgttgcgct catggtaatg gcagcactac acaattctct 540
taccgtcatg ccatccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtctgt 600
15ttgtgagtag tgtatacggc gaccaagctg ctcttgcccg gcgtctatac gggacaacac 660
cgcgccacat agcagtactt tgaaagtgc catcatcggg aatcgttctt cggggcggaa 720
agactcaagg atcttgccgc tattgagatc cagttcgata tagcccactc ttgcaccacg 780
ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
aaatgccgca aagaaggga tgagtgcgac acgaaaatgt tggatgctca tactcgtcct 900
20ttttcaatat tattgaagca tttatcaggg ttactagtac gtctctcaag gataagtaag 960
taatattaag gtacgggagg tattggacag gccgcaataa aatatcttta ttttcattac 1020
atctgtgtgt tgggtttttg tgtgaatcga tagtactaac atacgctctc catcaaaaca 1080
aaacgaaaca aaacaaacta gcaaaatagg ctgtccccag tgcaagtgca ggtgccagaa 1140
catttctctg gcctaactgg ccggtacc tg agctcgctag cctcgaggat atcaagatct 1200
25ggcctcggcg gccaaagctg gcaatccggg actgttggtg aagccaccat gg 1252
```

<210> 75

<400> 75

30 000

<210> 76

<211> 228

<212> DNA

35<213> Artificial Sequence

<220>

<223> A synthetic construct.

40<400> 76

```
actagtcgtc tctcttgaga gaccgcgac gccaccatga taagtaagta atattaaata 60
agtaaggcct gagtggccct cgagccagcc ttgagttggt tgagtccaag tcacgtctgg 120
aqaatctqta cctacqcqta aactctacat aactaaccqq ctcggcggcc gaattcttgc 180
```

45

gatctaagta agcttggcat tccggtactg ttggtaaagc caccatgg 228

<210> 77

<211> 228

5<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

10

<400> 77

actagtacgt ctctcttgag agaccgcat cgccaccatg ataagtaagt aatattaaat 60
aagtaaggcc tgagtggccc tcgagtccag ccttgagttg gt tgagtcca agtcacgtct 120
ggagatctgg taccttacgc gtagagctct acgtagctag cggcctcggc ggccgaattc 180
15ttgcgatcta agcttggcaa tccggtactg ttggtaaagc caccatgg 228

<210> 78

<211> 230

<212> DNA

20<213> Artificial Sequence

<220>

<223> A synthetic construct.

25<400> 78

actagtacgt ctctcttgag agaccgcat cgcattgccta ggtaggtagt attagagcat 60
aggtagaggc ctaagtggcc ctcgagtcca gccttgagtt ggtaggtcc aagtcacgtc 120
tggagatctg gtaccttacg cgtatgagct ctacgtagct agcggcctcg gcggccgaat 180
tcttgcatgc taagcttggc aatccggtac tgtaggtaaa gccaccatgg 230

30

<210> 79

<211> 234

<212> DNA

<213> Artificial Sequence

35

<220>

<223> A synthetic construct.

<400> 79

40actagtacgt ctctcttgag agaccgcat cgccaccatg tcttaggtagg tagtaaacga 60
aagggtctaa aggcctaagt ggccctcgag tccagccttg agttgggtga gtccaagtca 120
cgtttgagga tctgttacct tacgcgtatg agctctacgt agctagcggc ctcggcggcc 180
gaattcttgc gatctaagct tggcaatccg gtactgttgg taaagccacc atgg 234

46

<210> 80

<211> 938

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 80

```
10actagtaacc ctgataaatg cttcaataat attgaaaaag gaagagtatg agtat tcaac 60
   atttccgtgt cgcccttatt cccttttttg cggcattttg ccttcctggt tttgctcacc 120
   cagaaacgct ggtgaaagta aaagatgctg aagatcagtt gggcgcacga gtgggttaca 180
   tcgaactgga tctcaacagc ggtaagatcc ttgagagttt tcgccccgaa gaacgttttc 240
   caatgatgag cactttttaa gttctgctat gtggcgcggt attatcccggt attgacgccg 300
15ggcaagagca actcggctgc cgcatacact attctcagaa tgacttggtt gactactcac 360
   cagtcacaga aaagcatctt acggatggca tgacagtaag agaattatgc agtgctgcc 420
   taaccatgag tgataacacc gcggccaact tacttctgac aacgatcgga ggacccaagg 480
   agctaaccgc ttttttgac aacatggggg atcatgtaac tcgccttgat cgttggaac 540
   cggagctgaa tgaagccata ccaaacgacg agcgtgacac cacgatgcct gtagcaatgg 600
20caacaacggt gcgcaacta ttaactggcg aactacttac tctagcttcc cggcaacaat 660
   taatagactg gatggaggcg gataaagttg caggaccact tctgcgctcg gcccttccgg 720
   ctggctggtt tattgctgat aaatctggag ccggtgagcg tggctctgcg ggatcattg 780
   cagcactggg gccagatggt aagccctccc gtatcgtagt tatctacacg acggggagtc 840
   aggcaactat ggatgaacga aatagacaga tcgctgagat aggtgcctca ctgattaagc 900
25attggttaacc actgcagtgg ttttcctttt gcggccgc 938
```

<210> 81

<211> 938

<212> DNA

30<213> Artificial Sequence

<220>

<223> A synthetic construct.

35<400> 81

```
   actagtaacc ctgataaatg ctgcaaacat attgaaaaag gaagagtatg agtat tcaac 60
   atttccgtgt cgcactcatt cccttctttg cggcattttg cttgcctggt tttgcacacc 120
   ccgaaacgct ggtgaaagta aaagatgctg aagatcaact gggcgcacga gtgggttata 180
   tcgaactgga tctcaatagc ggtaagatcc ttgagagttt tcgccccgaa gaacgttttc 240
40caatgatgag cactttttaa gttctgctat gtggcgcggt attatcccggt attgacgccg 300
   ggcaagagca gctcggctgc cgcatacact actcacagaa cgacttggtt gactactcgc 360
   cggtcacgga aaagcatctt acggatggca tgacagtaag agaattgtgt agtgctgcc 420
   taaccatgag tgataacacc gcggccaact tacttctgac aacgatcgga ggccctaagg 480
```


47

```

agctgaccgc atttttgcac aacatggggg atcatgtaac ccggcttgat cgttgggaac 540
cggagctgaa cgaagccata ccgaacgacg agcgtgacac cacgatgcct gtagcaatgg 600
caacaacggt gcgcaaaact ctactggcg aacttctcac tctagcatca cgacagcaac 660
tcatagactg gatggaggcg gataaagttg caggaccact tctgcgctcg gcccttccgg 720
5ctggctggtt tatagctgat aaatccggtg ccggtgaacg cggctctcgc gggatcattg 780
ctgcgctggg gccagatggt aagccctcac gaatcgtagt tatctacacg acggggagtc 840
aggcaactat ggatgaacga aatagacaga tcgctgagat aggtgcctca ctgatcaagc 900
actggtagcc actgcagtgg tttagctttt gcggccgc 938

```

10<210> 82

<211> 938

<212> DNA

<213> Artificial Sequence

15<220>

<223> A synthetic construct.

<400> 82

```

actagtaacc ctgacaaatg ctgcaaacat attgaaaaag gaagagtatg agcatccaac 60
20attttcgtgt cgcactcatt cccttctttg cggcattttg cttgcctgtt tttgcacacc 120
ccgaaacgct ggtgaaagta aaagatgctg aagatcaact ggttgcaaga gtgggctata 180
tcgaactgga tctcaatagc ggcaagatcc ttgagtcttt tcgccccgaa gaacgttttc 240
cgatgatgag cactttttaa gttctgctat gtggcgcggt gttgtcccgt atagacgccg 300
ggcaagagca gcttggtcgc cgtatacact actcacaaaa cgacttggtt gagtactcgc 360
25cggctcacgga aaagcatctt acggatggca tgacggtaag agaattgtgt agtgctgcca 420
ttaccatgag cgacaatacc gcggccaact tacttctgac aacgatcgga ggccctaagg 480
agctgaccgc atttttgcac aacatggggg atcatgtaac ccggcttgac cgctgggaac 540
cggagctgaa cgaagccata ccgaacgacg agcgtgacac cacgatgcct gtagcaatgg 600
caacaacggt gcggaacta ctactggcg aacttctcac tctagcatca cgacagcagc 660
30tcatagactg gatggaggcg gacaaagtag caggaccact tcttcgctcg gccctccctg 720
ctggctggtt cattgctgat aaatccggtg ccggtgaacg cggctctcgc gggatcattg 780
ctgcgctggg gcctgatggt aagccctcac gaatcgtagt aatctacacg acggggagtc 840
aggccactat ggacgaacga aatagacaga tcgctgagat cgggtgcctca ctgatcaagc 900
actggttaacc actgcagtgg tttagcattt gcggccgc 938

```

35

<210> 83

<211> 938

<212> DNA

<213> Artificial Sequence

40

<220>

<223> A synthetic construct.

<400> 83

```
actagtaacc ctgacaaatg ctgcaaacat attgaaaaag gaagagtatg agcatccaac 60
atatttcgtgt cgcactcatt cccttccttg cggcattttg cttgcctgtt tttgcacacc 120
ccgaaaacgct ggtgaaagta aaagatgctg aagatcaact gggtgcaaga gtgggctata 180
5tcgaaactgga tctcaatagc ggcaagatcc ttgagtcttt ccgccccgaa gaacgttttc 240
cgatgatgag cactttcaaa gtactgctat gtggcgcggt gttgtccgt atagacgccg 300
ggcaagagca gcttggtcgc cgtatacact actcacaaaa cgacttggtt gactactcgc 360
cggtcacgga aaagcatctt acggatggca tgacggtaag agaattgtgt agtgctgcca 420
ttaccatgag cgataatacc gcggccaact tacttctgac aacgatcgga ggccctaagg 480
10agctgaccgc atttttgcac aacatgggtg atcatgtgac ccggcttgac cgctgggaac 540
cggagctgaa cgaagccata ccgaacgacg agcgtgacac cacgatgcct gtagcaatgg 600
caacaactct tcgaaacta ctactggcg aacttctcac tctagcatca cgacagcagc 660
tcatagactg gatggaggcg gacaaagtag caggaccact tcttcgctcg gccctccctg 720
ctggctggtt cattgctgat aaatctggag ccggtgagcg tggctctcgc ggtatcattg 780
15ctgcgctggg gcctgatggt aagccctcac gaatcgtagt aatctacacg acggggagtc 840
aggccactat ggacgaacga aatagacaga tcgctgagat cggtgccctca ctgatcaagc 900
actggttaacc actgcagtgg tttagcattt gcggccgc 938
```

<210> 84

20<211> 938

<212> DNA

<213> Artificial Sequence

<220>

25<223> A synthetic construct.

<400> 84

```
actagtaacc ctgacaaatg ctgcaaacat attgaaaaag gaagagtatg agcatccaac 60
atatttcgtgt cgcactcatt cccttccttg cggcattttg cttgcctgtt tttgcacacc 120
30ccgaaaacgct ggtgaaagta aaagatgctg aagatcaact gggtgcaaga gtgggctata 180
tcgaaactgga tctcaatagc ggcaagatcc ttgagtcttt ccgccccgaa gaacgattcc 240
cgatgatgag cactttcaaa gtactgctat gtggcgcggt gttgtccgt atagacgccg 300
ggcaagagca gcttggtcgc cgtatacact actcacaaaa cgacttggtt gactactcgc 360
cggtcacgga aaagcatctt acggatggca tgacggtaag agaattgtgt agtgctgcca 420
35ttaccatgag cgataatacc gcggccaact tacttctgac aacgatcgga ggccctaagg 480
agctgaccgc atttttgcac aacatgggtg atcatgtgac ccggcttgac cgctgggaac 540
cggagctgaa cgaagccata ccgaacgacg agcgtgatac cacgatgcca gtagcaatgg 600
ccacaactct tcgaaacta ctactggcg aacttctcac tctagcatca cgacagcagc 660
tcatagactg gatggaggcg gacaaagtag caggaccact tcttcgctcg gccctccctg 720
40ctggctggtt cattgctgac aaatccggtg ccggtgaacg cggctctcgc ggcattcattg 780
ctgcgctggg gcctgatggt aagccctcac gaatcgtagt aatctacacg acggggagtc 840
aggccactat ggacgaacga aatagacaga tcgctgagat cggtgccctca ctgatcaagc 900
actggttaacc actgcagtgg tttagcattt gcggccgc 938
```

<210> 85

<400> 85

000

5

<210> 86

<400> 86

000

10

<210> 87

<400> 87

000

15

<210> 88

<211> 1038

<212> DNA

<213> Artificial Sequence

20

<220>

<223> A synthetic construct.

<400> 88

25atgaagaagc ccgaactcac cgctaccagc gttgaaaaat ttctcatcga gaagttcgac 60
agtgtgagcg acctgatgca gttgtcggag ggcgaagaga gccgagcctt cagcttcgat 120
gtcggcgggac gcggtatgt actgcgggtg aatagctgcg ctgatggctt ctacaaagac 180
cgctacgtgt accgccactt cgccagcgct gcactacca tccccgaagt gttggacatc 240
ggcgagttca gcgagagcct gacatactgc atcagtagac gcgccaagg cgttactctc 300
30caagacctcc ccgaaacaga gctgcctgct gtgttacagc ctgtcgccga agctatggat 360
gctattgccg ccgccgacct cagtcaaacc agcggcttcg gcccatcgcg gccccaggc 420
atcggccagt acacaacctg gcgggatttc atttgcgcca ttgctgatcc ccatgtctac 480
cactggcaga ccgtgatgga cgacaccgtg tccgccagcg tagctcaagc cctggacgaa 540
ctgatgctgt gggccgaaga ctgtcccag gtgcgccacc tcgtccatgc cgacttcggc 600
35agcaacaacg tcctgaccga caacggccgc atcacgcgcg taatcgactg gtccgaagct 660
atgttcgggg acagtcagta cgaggtggcc aacatcttct tctggcggcc ctggctggct 720
tgcatggagc agcagactcg ctacttcgag cgccggcatc ccgagctggc cggcagccct 780
cgtctgcgag cctacatgct gcgcatcggc ctggatcagc tctaccagag cctcgtggac 840
ggcaacttcg acgatgctgc ctgggctcaa ggccgctgcg atgccatcgt ccgcagcggg 900
40gccggcaccg tcggtcgcac acaaatcgct cgccggagcg cagccgtatg gaccgacggc 960
tgctgcgagg tgctggccga cagcggcaac cgccggccca gtacacgacc gcgcgctaag 1020
gaggtaggtc gagtttaa 1038

<210> 89

<211> 4333

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 89

```
10ggcctaactg gccggtacct gagctcgcta gcctcgagga tatcaagatc tggcctcggc 60
   ggccaagctt ggcaatccgg tactgttggg aaagccacca tggaagatgc caaaaacatt 120
   aagaagggcc cagcgccatt ctaccactc gaagacggga ccgccggcga gcagctgcac 180
   aaagccatga agcgctacgc cctggtgccc ggcaccatcg cctttaccga cgcacatatc 240
   gaggtggaca ttacctacgc cgagtacttc gagatgagcg ttcggctggc agaagctatg 300
15aagcgctatg ggctgaatac aaaccatcgg atcgtggtgt gcagcgagaa tagcttgacg 360
   ttcttcattg ccgtgttggg tgccctgttc atcgtgtgtg ctgtggcccc agctaacgac 420
   atctacaacg agcgcgagct gctgaacagc atgggcatca gccagccac cgctcgattc 480
   gtgagcaaga aagggtgca aaagatcctc aacgtgcaaa agaagctacc gatcatacaa 540
   aagatcatca tcatggatag caagaccgac taccagggtt tccaaagcat gtacaccttc 600
20gtgacttccc atttgccacc cggcttcaac gactacgact tcgtgcccga gagcttcgac 660
   cgggacaaaa ccatcgccct gatcatgaac agtagtggca gtaccggatt gccaagggc 720
   gtagccctac cgcaccgac cgcttgtgtc cgattcagtc atgcccgcga ccccatcttc 780
   ggcaaccaga tcatccccga caccgctatc ctcagcgtgg tgccatttca ccacggcttc 840
   ggcattgttca ccacgtggg ctacttgatc tgcggctttc gggtcgtgct catgtaccgc 900
25ttcgaggagg agctattctt gcgcagcttg caagactata agattcaatc tgccctgctg 960
   gtgcccacac tathtagctt cttcgctaag agcactctca tcgacaagta cgacctaacg 1020
   aacttgacg agatcgccag cggcggggcg ccgctcagca aggaggtagg tgaggccgtg 1080
   gccaaacgct tccacctacc aggcattccg cagggctacg gcctgacaga aacaaccagc 1140
   gccattctga tcacccccga aggggacgac aagcctggcg cagtaggcaa ggtggtgccc 1200
30ttcttcgagg ctaagggtgg ggacttgac accggtaaga cactgggtgt gaaccagcgc 1260
   ggcgagctgt gcgtccgtgg ccccatgatc atgagcggct acgttaacaa ccccgaggct 1320
   acaaacgctc tcatcgacaa ggacggctgg ctgcacagcg gcgacatcgc ctactgggac 1380
   gaggacgagc acttcttcat cgtggaccgg ctgaagagcc tgatcaaata caagggctac 1440
   caggtagccc cagccgaact ggagagcatc ctgctgcaac accccaacat cttcgacgcc 1500
35ggggtcgccg gcctgcccga cgacgatgcc ggcgagctgc ccgccgcagt cgtcgtgctg 1560
   gaacacggta aaaccatgac cgagaaggag atcgtggact atgtggccag ccagggttaca 1620
   accgccaaag agctgcgcgg tgggtgtgtg ttcgtggacg aggtgcctaa aggactgacc 1680
   ggcaagttgg acgcccga gatccgcgag attctcatta aggccaaaga gggcggcaag 1740
   atcgccgtgt aataattcta gactcggggc ggccggccgc ttcgagcaga catgataaga 1800
40tacattgatg agtttgga aaccacaact agaatgcagt gaaaaaaatg ctttattttgt 1860
   gaaatttgtg atgctattgc tttatttgta accattataa gctgcaataa acaagttaac 1920
   aacaacaatt gcattcattt tatgtttcag gttcaggggg aggtgtggga ggttttttaa 1980
   agcaagtaaa acctctacaa atgtggtaaa atcgataagg atccgtcgac cgatgccctt 2040
```

gagagccttc aaccagtcg gtccttccg gtgggcgcgg ggcattgacta tcgtcgccgc 2100
 acttatgact gtcttcttta tcatgcaact cgtaggacag gtgcggcag cgctcttccg 2160
 ctctctcgct cactgactcg ctgcgctcgg tcgttcggct gcggcgagcg gtatcagctc 2220
 actcaaaggc ggtaatacgg ttatccacag aatcagggga taacgcagga aagaacatgt 2280
 5gagcaaaagg ccagcaaaag gccaggaacc gtaaaaaggc cgcgttgctg gcgtttttcc 2340
 ataggctccg cccccctgac gagcatcaca aaaatcgacg ctcaagtcag aggtggcgaa 2400
 acccgacagg actataaaga taccaggcgt ttccccctgg aagctccctc gtgcgctctc 2460
 ctgttccgac cctgccgctt accggatacc tgtccgcctt tctcccttcg ggaagcgtgg 2520
 cgctttctca tagctcacgc tgtaggatc tcagttcggg ttaggtcgtt cgctccaagc 2580
 10tggtgctgtgt gcacgaaccc cccgttcagc ccgaccgctg cgccttatcc ggtaactatc 2640
 gtcttgagtc caaccggta agacacgact tatcgccact ggcagcagcc actggtaaca 2700
 ggattagcag agcgaggtat gtaggcggtg ctacagagtt cttgaagtgg tggcctaact 2760
 acggctacac tagaagaaca gtatttggtg tctgcgctct gctgaagcca gttaccttcg 2820
 gaaaaagagt tggtagctct tgatccggca acaaaaccac cgctggtagc ggtggttttt 2880
 15ttgtttgcaa gcagcagatt acgcgcagaa aaaaaggatc tcaagaagat cctttgatct 2940
 tttctacggg gtctgacgct cagtggaaacg aaaactcacg ttaagggatt ttggtcatga 3000
 gattatcaaa aaggatcttc acctagatcc ttttaaatta aaaatgaagt tttaaatcaa 3060
 tctaaagtat atatgagtaa acttggctctg acagcggccg caaatgctaa accactgcag 3120
 tggttaccag tgcttgatca gtgaggcacc gatctcagcg atctgcctat ttcgttcgctc 3180
 20catagtggcc tgactccccg tcgtgtagat cactacgatt cgtgagggct taccatcagg 3240
 cccagcgcga gcaatgatgc cgcgagagcc gcgttcaccg gcccccgatt tgtcagcaat 3300
 gaaccagcca gcaggaggcg ccgagcgaag aagtggctct gctactttgt ccgcctccat 3360
 ccagtctatg agctgctgtc gtgatgctag agtaagaagt tcgccagtga gtagtttccg 3420
 aagagtgtg gccattgcta ctggcatcgt ggtatcacgc tcgtcgcttcg gtatggcttc 3480
 25gttcaactct gggtcccagc ggtcaagccg ggtcacatga tcaaccatat tatgaagaaa 3540
 tgcagtcagc tccttagggc ctccgatcgt tgtcagaagt aagttggccg cgggtgtgctc 3600
 gctcatggta atggcagcac tacacaattc tcttaccgtc atgccatccg taagatgctt 3660
 ttccgtgacc ggcgagtact caaccaagtc gttttgtgag tagtgtatac ggcgaccaag 3720
 ctgctcttgc ccggcgtcta tacgggacaa caccgcgcca catagcagta ctttgaaagt 3780
 30gctcatcacc gggaatcgtt cttcggggcg gaaagactca aggatcttgc cgctattgag 3840
 atccagttcg atatagccca ctcttgacc cagttgatct tcagcatctt ttactttcac 3900
 cagcgtttcg ggggtgtgcaa aaacaggcaa gcaaaatgcc gcaaagaagg gaatgagtgc 3960
 gacacgaaaa tgttgatgc tcatactcgt cctttttcaa tattattgaa gcatttatca 4020
 gggttactag tacgtctctc aaggataagt aagtaatat aaggtacggg aggtattgga 4080
 35caggccgcaa taaaatatct ttattttcat tacatctgtg tgttggtttt ttgtgtgaat 4140
 cgatagtact aacatacgtc ctccatcaaa acaaaacgaa acaaaacaaa ctagcaaaat 4200
 aggtgtccc cagtgcaggt gcaggcgcca gaacatttct ctaagtaata ttaaggtacg 4260
 ggaggtattg gacaggccgc aataaaatat ctttattttc attacatctg tgtgttggtt 4320
 ttttgtgtga atc 4333

<210> 90

<211> 3522

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 90

```
10ggcctaactg gccggtacct gagctcgcta gcctcgagga tatcaagatc tggcctcggc 60
   ggccaagctt ggcaatccgg tactgttggt aaagccacca tggcttccaa ggtgtacgac 120
   ccgagcaac gcaaacgcat gatcactggg cctcagtggg gggctcgctg caagcaaagt 180
   aacgtgctgg actccttcat caactactat gattccgaga agcacgccga gaacgccgtg 240
   atttttctgc atggtaacgc tgccctccagc tacctgtgga ggcacgtcgt gcctcacatc 300
15gagcccgtgg ctagatgcat catccctgat ctgatcggaa tgggtaagtc cggcaagagc 360
   gggaatggct catatcgctt cctggatcac tacaagtacc tcaccgcttg gttcgagctg 420
   ctgaaccttc caaagaaaat catctttgtg ggccacgact ggggggcttg tctggccttt 480
   cactactcct acgagcacca agacaagatc aaggccatcg tccatgctga gagtgcctg 540
   gacgtgatcg agtcctggga cgagtggcct gacatcgagg aggatatcgc cctgatcaag 600
20agcgaagagg gcgagaaaat ggtgcttgag aataaacttc tcgtcgagac catgtcccca 660
   agcaagatca tgcggaaact ggagcctgag gagttcgctg cctacctgga gccattcaag 720
   gagaagggcg aggttagacg gcctaccctc tcctggcctc gcgagatccc tctcgtaaag 780
   ggaggcaagc ccgacgtcgt ccgattgtc cgcaactaca acgcctacct tcgggccagc 840
   gacgatctgc ctaagatgtt catcgagtcc gaccctgggt tcttttccaa cgctattgtc 900
25gagggagcta agaagttccc taacaccgag ttcgtgaagg tgaagggcct ccacttcagc 960
   caggaggacg ctccagatga aatgggtaag tacatcaaga gcttcgtgga gcgcgtgctg 1020
   aagaacgagc agtaattcta gagtcggggc ggccggccgc ttcgagcaga catgataaga 1080
   tacattgatg agtttggaac aaccacaact agaatgcagt gaaaaaaatg ctttatttgt 1140
   gaaatttgtg atgctattgc tttatttgta accattataa gctgcaataa acaagttaac 1200
30aacaacaatt gcattcattt tatgtttcag gttcagggg aggtgtggga ggttttttaa 1260
   agcaagtaaa acctctacaa atgtggtaaa atcgataagg atccgtcgac cgatgccctt 1320
   gagagccttc aaccagtcga gtccttccg gtgggcgcgg ggcatgacta tcgtcgccgc 1380
   acttatgact gtcttcttta tcatgcaact cgtaggacag gtgccggcag cgctcttccg 1440
   cttcctcgct cactgactcg ctgcgctcgg tcgttcggct gcggcgagcg gtatcagctc 1500
35actcaaaggc ggtaatacgg ttatccacag aatcagggga taacgcagga aagaacatgt 1560
   gagcaaaaag ccagcaaaag gccaggaacc gtaaaaaggc cgcgttgctg gcgtttttcc 1620
   ataggctccg cccctcctgac gagcatcaca aaaatcgacg ctcaagttag aggtggcgaa 1680
   acccgacagg actataaaga taccaggcgt tccccctgg aagctccctc gtgcgctctc 1740
   ctgttccgac cctgccgctt accggatacc tgtccgcctt tctcccttcg ggaagcgtgg 1800
40cgcttttctc tagctcacgc tgtaggtatc tcagtccggg gtaggtcgtt cgctccaagc 1860
   tgggctgtgt gcacgaaccc ccggttcagc ccgaccgctg cgccttatcc ggtaactatc 1920
   gtcttgagtc caaccggta agacacgact tatcgccact ggcagcagcc actggtaaca 1980
   ggattagcag agcgaggtat gtaggcggtg ctacagagtt cttgaagtgg tggcctaact 2040
```

```

acggctacac tagaagaaca gtatttggtg tctgcgctct gctgaagcca gttaccttcg 2100
gaaaaagagt tggtagctct tgatccggca aacaaaccac cgctggtagc ggtgggttttt 2160
ttgtttgcaa gcagcagatt acgcgcagaa aaaaaggatc tcaagaagat cctttgatct 2220
tttctacggg gtctgacgct cagtggaaacg aaaactcacg ttaagggatt ttggtcatga 2280
5gattatcaaa aaggatcttc acctagatcc ttttaaatta aaaatgaagt tttaaatcaa 2340
tctaaagtat atatgagtaa acttggctctg acagcggccg caaatgctaa accactgcag 2400
tggttaccag tgcttgatca gtgaggcacc gatctcagcg atctgcctat ttcgttcgtc 2460
catagtggcc tgactccccg tcgtgtagat cactacgatt cgtgagggct taccatcagg 2520
ccccagcgca gcaatgatgc cgcgagagcc gcgttcaccg gcccccgatt tgtcagcaat 2580
10gaaccagcca gcagggaggg ccgagcgaag aagtggctct gctactttgt ccgcctccat 2640
ccagtctatg agctgctgtc gtgatgctag agtaagaagt tcgccagtga gtagtttccg 2700
aagagttgtg gccattgcta ctggcatcgt ggtatcacgc tcgtcgttcg gtatggcttc 2760
gttcaactct ggttcccagc ggtcaagccg ggtcacatga tcacccatat tatgaagaaa 2820
tgcagtcagc tccttagggc ctccgatcgt tgtcagaagt aagttggccg cgggtgtgtc 2880
15gctcatggta atggcagcac tacacaattc tcttaccgtc atgccatccg taagatgctt 2940
ttccgtgacc ggcgagtact caaccaagtc gttttgtgag tagtgatac ggcgaccaag 3000
ctgctcttgc ccggcgtcta tacgggacaa caccgcgcca catagcagta ctttgaaagt 3060
gctcatcatc gggaatcgtt cttcggggcg gaaagactca aggatcttgc cgctattgag 3120
atccagttcg atatagccca ctcttgccacc cagttgatct tcagcatctt ttactttcac 3180
20cagcgtttcg ggggtgtgcaa aaacaggcaa gcaaaatgcc gcaaagaagg gaatgagtgc 3240
gacacgaaaa tgttggtatgc tcatactcgt cttttttcaa tattattgaa gcatttatca 3300
gggttactag tacgtctctc aaggataagt aagtaatatt aaggtagcgg aggtattgga 3360
caggccgcaa taaaatatct ttattttcat tacatctgtg tgttggtttt ttgtgtgaat 3420
cgatagtact aacatacgct ctccatcaaa acaaaacgaa acaaaacaaa ctagcaaaat 3480
25aggctgtccc cagtgcgaagt gcaggtgcca gaacatttct ct 3522

```

<210> 91

<211> 621

<212> DNA

30<213> Artificial Sequence

<220>

<223> A synthetic construct.

35<400> 91

```

gctagcgcca ccatgaccga gtacaagccc accgtgcgcc tggccaccgc cgacgacgtg 60
ccccgcgcgc tgccgaccct ggccgcgcgc ttcgcgcgact accccgccac ccgccacacc 120
gtggaccgcc accgccacat cgagcgcgtg accgagctgc aggagctggt cctgaccgc 180
gtgggccttg acatcggcaa ggtgtgggtg gccgacgacg gcgcgcgcgt ggccgtgtgg 240
40accacccccg agagcgtgga ggccggcgcc gtgttcgccg agatcgcccc ccgcatggcc 300
gagctgagcg gcagccgcct ggccgcccag cagcagatgg agggcctgct ggccccccac 360
cgccccaaag agcccgcctg gttcctggcc accgtgggcg tgagccccga ccaccagggc 420
aagggccttg gcagcgccgt ggtgctgccc ggcgtggagg ccgccgagcg cgccggcgtg 480

```

54

cccgccttcc tggagaccag cgtccccgc aacctgccct tctacgagcg cctgggcttc 540
acctgaccg ccgacgtgga ggtgcccag ggtccccga cctggtgcat gacctgcaag 600
cccgccgcct aatgatctag a 621

5<210> 92

<211> 621

<212> DNA

<213> Artificial Sequence

10

<220>

<223> A synthetic construct.

<400> 92

15gctagcgcca ccatgaccga gtacaagcct acctgcgcc tggccactcg cgatgatgtg 60
ccccgcgccg tccgactct ggccgccgt ttcgccgact acccgctac ccggcacacc 120
gtggaccccg accggcacat cgagcgtgtg acagagttgc aggagctgtt cctgacccgc 180
gtcgggcttg acatcggaac ggtgtgggta gccgacgacg gcgcgccgt ggccgtgttg 240
actaccccg agagcgttga ggccggcgcc gtgttcgccg agatcgcccc ccgaatggcc 300
20gagctgagcg gcagccgcct ggccgccag cagcaaatgg agggcctgct tgccccccat 360
cgtcccaagg agccgcctg gtttctggcc actgtaggag tgagccccga ccaccagggc 420
aagggcttg gcagcgccgt cgtgttgccc ggcgtagagg ccgccgaacg cgccggtgtg 480
ccgccttctc tggagacaag cgctccgct aaccttccat tctacgagcg cctgggcttc 540
acctgaccg ccgatgtcga ggtgcccag ggacccgga cctggtgcat gactcgcaag 600
25cctggcgccct aatgatctag a 621

<210> 93

<211> 621

<212> DNA

30<213> Artificial Sequence

<220>

<223> A synthetic construct.

35<400> 93

gctagcgcca ccatgaccga gtacaagcct acctgcgcc tggccactcg cgatgatgtg 60
ccccgcgccg tccgactct ggccgccgt ttcgccgact acccgctac ccggcacacc 120
gtggaccccg accggcacat cgagcgtgtg acagagttgc aggagctgtt cctgacccgc 180
gtcgggcttg acatcggaac ggtgtgggta gccgacgacg gcgcgccgt ggccgtgttg 240
40actaccccg agagcgttga ggccggcgcc gtgttcgccg agatcgcccc ccgaatggcc 300
gagctgagcg gcagccgcct ggccgccag cagcaaatgg agggcctgct tgccccccat 360
cgtcccaagg agcctgcctg gtttctggcc actgtaggag tgagccccga ccaccagggc 420
aagggcttg gcagcgccgt cgtgttgccc ggcgtagagg ccgccgaacg cgccggtgtg 480

55

```
cccgccctttc tcgaaacaag cgcaccaaga aaccttccat tctacgagcg cctgggcttc 540
accgtgaccg ccgatgtcga ggtgcccagag ggacctagga cctggtgtat gacacgaaaa 600
cctggcgcct aatgatctag a 621
```

5<210> 94

<211> 1672

<212> DNA

<213> Artificial Sequence

10<220>

<223> A synthetic construct.

<400> 94

```
aaagccacca tggaagatgc caaaaacatt aagaaggggc ctgctccctt ctaccctctt 60
15gaagatggga ctgctggcga gcaacttcac aaagctatga agcggatgac tcttgtgcca 120
gggacaattg cggttcacgga tgctcacatt gaagtagaca tcacatacgc tgagtatttt 180
gagatgtcgg tgccggtggc agaagctatg aagcgctatg ggctgaatac aaaccataga 240
attgtagtgt gcagtggaaa ctcggtgcag ttctttatgc ccgtgctggg ggctctcttc 300
atcgggggtg ctgtggctcc tgctaacgac atctacaacg agcgagagct gttgaactcg 360
20atggggatct ctcagcctac agtgggtgtt gtgagtaaga aagggtctca aaagattctc 420
aatgtgcaaa agaagctgcc tattatacaa aagattatta ttatggactc taagacagac 480
taccaggggt ttcagtcctc gtacacattt gtaacctctc atctgcctcc tggcttcaac 540
gagtacgact tcgtgcccga gtctttcgac agggacaaaa cgattgctct gatcatgaac 600
agctccgggt ctaccgggtc gcctaagggt gtagctctgc cccatcgaa agcttgtgtg 660
25agattctctc atgccaggga cccgatcttt ggaaaccaga tcacccctga cactgctatt 720
ctgtcgggtg tgccctttca tcatgggttt gggatgttca caacactggg atacctcatt 780
tgccgggtta gagtgggtgt catgtatagg tttgaagaag aactattcct acgctctttg 840
caagattata agattcagtc tgctctgctg gtgccaacac tattctcttt ttttgctaag 900
tctacgctca tagacaagta tgacttgctc aacttgacg agattgcttc tggcggagca 960
30cctctgtcta aggaggtagg tgaggctgtg gctaagcgtc ttcattctgc tgggtatcaga 1020
caggggtacg ggctaacaga acaacttct gctattctga ttacaccaga gggcgatgac 1080
aaaccgggg ctgtagggaa agtggtgccc ttttttgaag ccaaagtagt tgatcttgat 1140
accggtgaaga cactaggggt gaaccagcgt ggtgaactgt gtgtgcgggg ccctatgatt 1200
atgtcgggg acgttaacaa cccgaagct acaaatgctc tcatagacaa ggacgggtg 1260
35cttcatagcg gcgacattgc ctactgggac gaggatgagc atttcttcat cgtggacaga 1320
ctgaagtcgt tgatcaaata caaggggtat caagtagctc ctgccgagct tgagtccatt 1380
ctgcttcaac accccaatat ctctgatgct ggggtggctg ggctgcctga tgatgatgct 1440
ggagagctgc ctgctgctgt agtagtgctt gagcatggta agacaatgac agagaaggag 1500
atcgtggatt atgtggcttc acaagtgaac acagctaaga aactccgagg tggcgttgtg 1560
40tttgtggatg aggtgcctaa agggctcact ggcaagctgg atgccagaaa aattcgagag 1620
attctcatta aggctaagaa ggggtggaag attgctgtgt aatagttcta ga 1672
```

56

<210> 95

<211> 1166

<212> DNA

<213> Artificial Sequence

5

<220>

<223> A synthetic construct.

<400> 95

```
10gcggccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
   ctcagcgatc tgtctatttc gttcgtccat agtggcctga ctccccgtcg tgtagattac 120
   tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
   ttcaccggca cgggatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
   tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
15gagaagttcg ccagtgahta gtttcgaag agttgtggcc attgctactg gcatcggtgt 360
   atcacgctcg tcgttcggta tggcttcgtt cagctccggt tcccagcggg caagccgggt 420
   cacatgatca cccatgttgt gcaaaaatgc ggtcagctcc ttagggcctc cgatcggtgt 480
   cagaagtaag ttggccgcgg tattatcgct catggtaatg gcagcactac acaattctct 540
   taccgtcatg ccaccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtcgtt 600
20ttgtgagtag tgtatacggc gaccaagctg ctcttgcccg gcgtctatac gggacaacac 660
   cgcgccacat agcagtactt tgaaagtgt catcatcggg aatcggttctt cggggcggaa 720
   agactcaagg atcttgccgc tattgagatc cagttcgata tagccactc ttgcaccag 780
   ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
   aaatgcgcga aagaaggga tgagtgcgac acgaaaatgt tggatgctca tactcttct 900
25ttttcaatat gtttcgagca tttgtcaggg ttactagtac gtctctcttg agagaccgcg 960
   atcgccacca tgtctaggta ggtagtaaac gaaagggtt aaaggcctaa gtggccctcg 1020
   agtcagcct tgagttggtt gagtccaagt cacgtttgga gatctggtac cttacgcgta 1080
   tgagctctac gtagctagcg gcctcggcgg ccgaattctt gcgatctaag cttggcaatc 1140
   cggtactgtt ggtaaagcca ccatgg                                     1166
```

30

<210> 96

<211> 1166

<212> DNA

<213> Artificial Sequence

35

<220>

<223> A synthetic construct.

<400> 96

```
40gcggccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
   ctcagcgatc tgtctatttc gttcgtccat agtggcctga ctccccgtcg tgtagattac 120
   tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
   ttcaccggcc cccgatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
```

57

```

tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
aagaagttcg ccagttagta gtttccgaag agttgtggcc attgctactg gcatcgtggg 360
atcacgctcg tcgttcggta tggcttcggt caactccggt tcccagcggg caagccgggt 420
cacatgatca cccatgttgt gcaaaaatgc ggtcagctcc ttagggcctc cgatcgttgt 480
5cagaagtaag ttggccgcgg tgttgtcgct catggtaatg gcagcactac acaattctct 540
taccgtcatg ccatccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtcgtt 600
ttgtgagtag tgtatacggc gaccaagctg ctcttgcccg gcgtctatac gggacaacac 660
cgcgccacat agcagtactt tgaaagtgt catcatcggg aatcgttctt cggggcggaa 720
agactcaagg atcttgccgc tattgagatc cagttcgata tagccactc ttgcacccag 780
10ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
aaatgccgca aagaagggaa tgagtgcgac acgaaaatgt tggatgctca tactcttcct 900
ttttcaatat gtttgcagca tttgtcaggg ttactagtac gtctctcttg agagaccgcg 960
atcgccacca tgtctaggta ggtagtaaac gaaagggcct aaaggcctaa gtggccctcg 1020
agtccagcct tgagttgggt gagtccaagt cacgtttgga gatctggtac cttacgcgta 1080
15tgagctctac gtagctagcg gcctcggcgg ccgaattctt gcgttcgaag cttggcaatc 1140
cgttactgtt ggtaaagcca ccatgg 1166

```

<210> 97

<211> 1166

20<212> DNA

<213> Artificial Sequence

<220>

<223> A synthetic construct.

25

<400> 97

```

gcgccgcaa atgctaaacc actgcagtgg ttaccagtgc ttgatcagtg aggcaccgat 60
ctcagcgatc tgcctatttc gttcgtccat agtggcctga ctccccgtcg tgtagatcac 120
tacgattcgt gagggcttac catcaggccc cagcgcagca atgatgccgc gagagccgcg 180
30ttcaccggcc cccgatttgt cagcaatgaa ccagccagca gggagggccg agcgaagaag 240
tggtcctgct actttgtccg cctccatcca gtctatgagc tgctgtcgtg atgctagagt 300
aagaagttcg ccagttagta gtttccgaag agttgtggcc attgctactg gcatcgtggg 360
atcacgctcg tcgttcggta tggcttcggt caactctggt tcccagcggg caagccgggt 420
cacatgatca cccatgttgt gcaaaaatgc ggtcagctcc ttagggcctc cgatcgttgt 480
35cagaagtaag ttggccgcgg tgttgtcgct catggtaatg gcagcactac acaattctct 540
taccgtcatg ccatccgtaa gatgcttttc cgtgaccggc gagtactcaa ccaagtcgtt 600
ttgtgagtag tgtatacggc gaccaagctg ctcttgcccg gcgtctatac gggacaacac 660
cgcgccacat agcagtactt tgaaagtgt catcatcggg aatcgttctt cggggcggaa 720
agactcaagg atcttgccgc tattgagatc cagttcgata tagccactc ttgcacccag 780
40ttgatcttca gcatctttta ctttcaccag cgtttcgggg tgtgcaaaaa caggcaagca 840
aaatgccgca aagaagggaa tgagtgcgac acgaaaatgt tggatgctca tactcttcct 900
ttttcaatat gtttgcagca tttgtcaggg ttactagtac gtctctcttg agagaccgcg 960
atcgccacca tgtctaggta ggtagtaaac gaaagggcct aaaggcctaa gtggccctcg 1020

```

58

agtccagcct tgagttgggt gagtccaagt cacgtttgga gatctggtac cttacgcgta 1080
tgagctctac gtagctagcg gcctcggcgg ccgaattctt gcgttcgaag cttggcaatc 1140
cgg tactggt ggtaaagcca ccatgg 1166

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
30 March 2006 (30.03.2006)

PCT

(10) International Publication Number
WO 2006/034061 A3

(51) International Patent Classification:
C12N 15/09 (2006.01) C12N 15/31 (2006.01)

(21) International Application Number:
PCT/US2005/033218

(22) International Filing Date:
16 September 2005 (16.09.2005)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
10/943,508 17 September 2004 (17.09.2004) US

(71) Applicant (for all designated States except US):
PROMEGA CORPORATION [US/US]; 2800 Woods
Hollow Road, Madison, WI 53711 (US).

(72) Inventors; and

(75) Inventors/Applicants (for US only): **WOOD, Keith, V.**
[US/US]; 8380 Swan Road, Mt. Horeb, WI 53572 (US).
WOOD, Monika, G. [US/US]; 8380 Swan Road, Mt.
Horeb, WI 53572 (US). **ALMOND, Biran** [US/US]; 5765
Richard Drive, Fitchburg, WI 53719 (US). **PAGUIO,**
Aileen [US/US]; 205 Ramsey Court, Madison, WI 53704
(US). **FAN, Frank** [TZ/US]; 2977 Dunmore Street, Madi-
son, WI 53711 (US).

(74) Agents: **STEFFEY, Charles, E.** et al.; Schwegman, Lund-
berg, Woessner & Kluth, P.A., P.O. Box 2938, Minneapolis,
MN 55402 (US).

(81) Designated States (unless otherwise indicated, for every
kind of national protection available): AE, AG, AL, AM,
AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN,
CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI,
GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE,
KG, KM, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, LY,
MA, MD, MG, MK, MN, MW, MX, MZ, NA, NG, NI, NO,
NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK,
SL, SM, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ,
VC, VN, YU, ZA, ZM, ZW.

(84) Designated States (unless otherwise indicated, for every
kind of regional protection available): ARIPO (BW, GH,
GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM,
ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM),
European (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI,
FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, NL, PL, PT,
RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA,
GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

- with international search report
- before the expiration of the time limit for amending the
claims and to be republished in the event of receipt of
amendments

(88) Date of publication of the international search report:
26 May 2006

For two-letter codes and other abbreviations, refer to the "Guid-
ance Notes on Codes and Abbreviations" appearing at the begin-
ning of each regular issue of the PCT Gazette.



WO 2006/034061 A3

(54) Title: SYNTHETIC NUCLEIC ACID MOLECULE AND METHODS OF PREPARATION

(57) Abstract: A method to prepare synthetic nucleic acid molecules having reduced inappropriate or unintended transcriptional characteristics when expressed in a particular host cell.

INTERNATIONAL SEARCH REPORT

 International application No
 PCT/US2005/033218

A. CLASSIFICATION OF SUBJECT MATTER

C12N15/09 C12N15/31

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

C12N

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, BIOSIS, EMBASE, Sequence Search

C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|-----------|---|-----------------------|
| X | WO 01/23541 A (ALEXION PHARMACEUTICALS, INC; FODOR, WILLIAM, L; RAMSOONDAR, JAGDEECE,) 5 April 2001 (2001-04-05) * There is 94,521% identity in 803 nt overlap (total 825 nt) of the sequence shown in Fig. 4 with SEQ ID NO: 4 of the present application * | 11,15 |
| A | WO 2004/042010 A (UNIVERSITY OF TENNESSEE RESEARCH FOUNDATION) 21 May 2004 (2004-05-21) | |
| A | US 5 670 356 A (SHERF ET AL) 23 September 1997 (1997-09-23) cited in the application | |

-/--

☒ Further documents are listed in the continuation of Box C.☒ See patent family annex.

* Special categories of cited documents :

- *A* document defining the general state of the art which is not considered to be of particular relevance
- *E* earlier document but published on or after the International filing date
- *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- *O* document referring to an oral disclosure, use, exhibition or other means
- *P* document published prior to the International filing date but later than the priority date claimed

- *T* later document published after the International filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.
- *Z* document member of the same patent family

Date of the actual completion of the International search

9 March 2006

Date of mailing of the International search report

31.03.06

Name and mailing address of the ISA/

 European Patent Office, P.B. 5818 Patentlaan 2
 NL - 2280 HV Rijswijk
 Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
 Fax: (+31-70) 340-3016

Authorized officer

Hillenbrand, G

INTERNATIONAL SEARCH REPORT

 International application No
 PCT/US2005/033218

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|-----------|---|-----------------------|
| X | US 2002/100076 A1 (GARCON FREDERIC ET AL) 25 July 2002 (2002-07-25) * There is 88,19% identity of SEQ ID NO: 2 in 906 nt overlap (total 5909 nt) with SEQ ID NO: 74 (1252 nt) of the present application - 39,2% identity of SEQ ID NO: 2 to SEQ ID NO: 41 * | 47,49,50 |
| X | WO 97/08320 A (MORPHOSYS GESELLSCHAFT FUER PROTEINOPTIMIERUNG MBH; KNAPPIK, ACHIM; PA) 6 March 1997 (1997-03-06) * There is 88,19% identity of the sequence of Fig. 36 in 906 nt overlap (total 1289 nt) with SEQ ID NO: 74 (1252 nt) of the present application - 39,2% identity of SEQ ID NO: 2 to SEQ ID NO: 41 * | 47,49,50 |
| X | DATABASE EMBL 1 March 1996 (1996-03-01), GROSKREUTZ ET AL.: "Cloning vector pGL3-Basic, complete sequence" XP002371236 retrieved from EBI Database accession no. U47295 * There is 85,82% identity of U47295 in 3095 nt overlap (total 4818 nt) with SEQ ID NO: 89 (4333 nt) of the present application * abstract | 63-67 |
| X | DATABASE EMBL 15 May 2001 (2001-05-15), ZHUANG, Y. ET AL.: "Co-reporter vector phRG-B, complete sequence" XP002371237 retrieved from EBI Database accession no. AF362550 * There is 98,82% identity of AF362550 in 2375 nt overlap (total 4101 nt) with SEQ ID NO: 90 (3522 nt) of the present application * abstract | 63-67 |

INTERNATIONAL SEARCH REPORT

International application No.
CT/US2005/033218

Box II Observations where certain claims were found unsearchable (Continuation of item 2 of first sheet)

This International Search Report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:
because they relate to subject matter not required to be searched by this Authority, namely:
2. ☒ Claims Nos.: 1-10, 12-14, 16-30, 32-46, 48, 53, 55-62, 68-69
because they relate to parts of the International Application that do not comply with the prescribed requirements to such an extent that no meaningful International Search can be carried out, specifically:
see FURTHER INFORMATION sheet PCT/ISA/210
3. ☐ Claims Nos.:
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box III Observations where unity of invention is lacking (Continuation of item 3 of first sheet)

This International Searching Authority found multiple inventions in this International application, as follows:

see additional sheet

1. ☐ As all required additional search fees were timely paid by the applicant, this International Search Report covers all searchable claims.
2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.
3. ☒ As only some of the required additional search fees were timely paid by the applicant, this International Search Report covers only those claims for which fees were paid, specifically claims Nos.:
11 and 15 (partially), 47, 49 and 50 (partially), 63-67 (partially)
4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this International Search Report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☐ The additional search fees were accompanied by the applicant's protest.
- ☒ No protest accompanied the payment of additional search fees.

FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

Continuation of Box II.2

Claims Nos.: 1-10, 12-14, 16-30, 32-46, 48, 53, 55-62, 68-69

The present application contains 69 claims, of which 7 claims are independent. They are drafted in such a way that the claims as a whole are not in compliance with the provisions of clarity and conciseness of Article 6 PCT, as they erect a smoke screen in front of the skilled reader when assessing the intended scope of protection. In view of the fact that the starting (parent) nucleic acid sequences are not defined in most claims, it is impossible for the skilled reader to determine the subject-matter for which protection is sought. The non-compliance with the substantive provisions of the PCT is to such an extent, that a meaningful search of the claims identified above was not possible.

The applicant's attention is drawn to the fact that claims relating to inventions in respect of which no international search report has been established need not be the subject of an international preliminary examination (Rule 66.1(e) PCT). The applicant is advised that the EPO policy when acting as an International Preliminary Examining Authority is normally not to carry out a preliminary examination on matter which has not been searched. This is the case irrespective of whether or not the claims are amended following receipt of the search report or during any Chapter II procedure. If the application proceeds into the regional phase before the EPO, the applicant is reminded that a search may be carried out during examination before the EPO (see EPO Guideline C-VI, 8.5), should the problems which led to the Article 17(2) declaration be overcome.

FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

This International Searching Authority found multiple (groups of) inventions in this international application, as follows:

Inventions 1-20 : claims 11 and 15 (partially)

The subject-matter of this group of different inventions comprises an isolated nucleic acid molecule comprising a synthetic nucleotide sequence having a coding region for a selectable polypeptide, wherein the synthetic nucleotide sequence has 90% or less nucleic acid sequence identity to a parent nucleotide encoding a corresponding selectable polypeptide, wherein the nucleotide sequence encodes a selectable polypeptide with at least 85% amino acid sequence identity to the corresponding selectable polypeptide encoded by the the parent nucleotide sequence - wherein the synthetic nucleotide sequence comprises an open reading frame in SEQ ID NO: 4 to SEQ ID NO: 84 as claimed in claims 11 and 15.

Invention 21: claim 31 (partially)

The subject-matter of this invention comprises an isolated nucleic acid sequence encoding a firefly luciferase, wherein the synthetic nucleotide sequence has 80% or less nucleic acid sequence identity to a parent nucleotide having SEQ ID NO: 43 or 85% or less nucleic acid sequence identity to a parent nucleic acid sequence having SEQ ID NO: 14 which encodes a firefly luciferase, wherein the nucleotide sequence encodes a firefly luciferase with at least 85% amino acid sequence identity to the corresponding luciferase encoded by the the parent nucleotide sequence, wherein the synthetic nucleotide sequence comprises an open reading frame in SEQ ID NO: 21-23.

Invention 23: claims 47, 49 and 50 (partially)

A plasmid comprising SEQ ID NO: 74 which comprises an open reading frame with less than 90% nucleic acid sequence identity to 41 which confers resistance to ampicillin.

Inventions 24-46: claims 51-52 (partially)

A polynucleotide which hybridizes under stringent hybridization conditions to SEQ ID NO: 4 to SEQ ID NO: 23 as claimed in claim 51 and encodes a selectable polypeptide or a firefly luciferase.

Invention 47: claim 54

FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

An isolated nucleic acid molecule comprising a synthetic nucleotide sequence which does not code for a desirable peptide or polypeptide but includes sequences which inhibit transcription and/or translation wherein the synthetic nucleotide sequence has SEQ ID NO: 49.

Inventions 48-49: claims 63-67 (partially)

A plasmid which includes a sequence including SEQ ID NO : 89 or SEQ ID NO: 90.
The search was limited to matter related to invention 1 and inventions 23, 48 and 49 as requested by the applicant in his letter dated 13.02.2006.

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No
PCT/US2005/033218

| Patent document cited in search report | | Publication date | Patent family member(s) | Publication date |
|---|----|---------------------|----------------------------|---------------------|
| WO 0123541 | A | 05-04-2001 | AU 7744800 A | 30-04-2001 |
| | | | CA 2385162 A1 | 05-04-2001 |
| | | | EP 1220928 A2 | 10-07-2002 |
| | | | JP 2003510072 T | 18-03-2003 |
| | | | MX PA02003232 A | 22-09-2003 |
| WO 2004042010 | A | 21-05-2004 | AU 2003301883 A1 | 07-06-2004 |
| US 5670356 | A | 23-09-1997 | NONE | |
| US 2002100076 | A1 | 25-07-2002 | AT 306553 T | 15-10-2005 |
| | | | BR 0104564 A | 04-06-2002 |
| | | | EP 1186666 A2 | 13-03-2002 |
| | | | FR 2812883 A1 | 15-02-2002 |
| WO 9708320 | A | 06-03-1997 | AT 219517 T | 15-07-2002 |
| | | | AU 725609 B2 | 12-10-2000 |
| | | | AU 6874596 A | 19-03-1997 |
| | | | CA 2229043 A1 | 06-03-1997 |
| | | | DE 69621940 D1 | 25-07-2002 |
| | | | DE 69621940 T2 | 16-01-2003 |
| | | | DK 859841 T3 | 09-09-2002 |
| | | | ES 2176484 T3 | 01-12-2002 |
| | | | JP 2001519643 T | 23-10-2001 |
| | | | PT 859841 T | 29-11-2002 |
| | | | US 6300064 B1 | 09-10-2001 |